

This document summarizes Siyuan Tang’s exposure to and understanding of research areas of interest. The application system does not provide a section to submit such a document.

Understanding of Research

Siyuan Tang

Last modified on February 1, 2025.

“All models are wrong, but some are useful.” ([Box and Draper \(1987\)](#))

Contents

1	Research Internship at the University of Michigan	2
1.1	Research Experience in Ann Arbor	2
1.2	Working Remotely After Returning to China	3
1.2.1	Analyzing Flare Dependency Structure via Bayesian Hierarchical HMMs	3
1.3	Lessons from the Internship	7
2	Statistical Network Analysis	7
2.1	The Stochastic Block Model	7
2.1.1	Why Estimation via EM is Intractable	8
2.2	A Pseudo-Likelihood Method	9
2.2.1	Estimation via Alternating Update	9
2.3	Jin’s SCORE Algorithm	11
2.3.1	The Oracle Case	12
2.3.2	The Real Case	14
2.4	Extensions of the SBM	14
2.5	Prediction Models for Network-Linked Data	15
2.5.1	Linear Regression with Network Cohesion	16
2.5.2	A Bayesian Interpretation	17
2.5.3	Prediction and Choosing the Tuning Parameter	18
3	Statistical Text Analysis	18
3.1	Mixture of Unigrams	19
3.1.1	Estimation via EM	19
3.2	Hofmann’s pLSI Model	21
3.2.1	An Equivalent Formulation	22
3.3	Latent Dirichlet Allocation	23
3.3.1	Estimation via Variational Inference	25

4	Covariance Estimation	26
4.1	A Linear Shrinkage Estimator	27
4.2	Generalized Thresholding	29
4.2.1	Proof of Consistency	29
4.3	The POET Estimator	32
5	Deep Learning	33
5.1	Diffusion Models	33
5.1.1	DDPM	34
5.1.2	A Score-Based Interpretation	36
5.1.3	Theoretical Foundation	37
5.2	Large Language Models	38
6	Statistics in the New Era	39

1 Research Internship at the University of Michigan

1.1 Research Experience in Ann Arbor

I was an on-site research intern working in Professor [Yang Chen](#)'s group at the Department of Statistics, University of Michigan (UM), Ann Arbor, Michigan, from July 2024 to November 2024.

Our work focused on predicting the energy released by solar flares and utilizing this information to classify flares into pre-defined categories. Our focus extended beyond regression performance to include binary classification performance, measured by the True Skill Statistic (TSS; [Woodcock \(1976\)](#)). Early on in this project I encountered a challenge: some solar flares had ill-defined and extremely long durations, complicating the analysis of their dynamics. To address this, I applied a truncation method, focusing on the critical moments surrounding the peak times of the flares, enabling more accurate modeling of their dynamics. I then proposed a profile modeling strategy to effectively denoise the observed sequences by employing a trend function that assumed linear growth before the peak and linear decay after. I fitted the model and conducted thorough goodness-of-fit assessments by evaluating residual sequences and testing for independence and normality.

Additionally, I explored how energy release patterns could enhance the classification of flare events, especially for stronger ones. Our work marked the first statistical analyses of solar flares from an energy perspective and illustrated that energy-based classifications often outperform traditional intensity-based methods ([Jiao et al. \(2020\)](#)). Currently, the team is preparing [a paper that summarizes our main findings](#).

1.2 Working Remotely After Returning to China

I continued to work after coming back to China in November 2024. Part of my research focused on developing robust binary classification methods that perform well under distribution shift, particularly for solar flare prediction tasks. I explored two types of distribution shift:

- Label shift, where the proportion of positive cases varies between training and test data.
- Covariate shift, where the feature distribution changes.

For label shift, I proved that TSS optimization is inherently invariant to changes in class proportions, and derived the closed form of the optimal classifier that maximizes TSS. To overcome the discontinuity of TSS, I developed a smooth approximation using sigmoid functions. My theoretical results were verified on synthetic data. Through real data experiments on the GOES event list ([Garcia \(1994\)](#)), I demonstrated that TSS optimization maintains robust performance under varying degrees of label shift, often outperforming traditional classifiers like Support Vector Machines ([Cortes \(1995\)](#)), Linear Discriminant Analysis, and Quadratic Discriminant Analysis.

For covariate shift, I implemented an importance weighting strategy to correct for the distribution difference between training and test data, deriving a modified TSS objective that incorporates these weights. This methodology was inspired by the textbook [Sugiyama and Kawanabe \(2012\)](#) and recent work from Professor [Tianxi Cai](#)'s group [Wang et al. \(2024\)](#). However, the results were not better than those for label shift, suggesting that the label shift assumption may be more appropriate for solar flare prediction.

1.2.1 Analyzing Flare Dependency Structure via Bayesian Hierarchical HMMs

In addition to the aforementioned work, I also sketched a hierarchical hidden markov model (HMM; [Rabiner and Juang \(1986\)](#)) to account for heterogeneity of various phases of flaring and various flaring mechanisms, whose details will be discussed in the sequel.

We consider a 6-hour non-overlapping floating window for each trajectory, and identify $M = 3$ kinds of observations:

- Observation 3: M or X flares.
- Observation 2: C flares.
- Observation 1: no flares or B flares.

See Table 1 for different flare types and their potential effects on Earth.

Let K denote the number of latent states. We assume that each trajectory is a realization of an HMM, with a shared observation matrix and initial state distribution across all trajectories, while allowing sample-specific variations in the state transition matrix. Specifically, for a trajectory indexed by l ($1 \leq l \leq N$, where N is the total number of active regions), its individual parameters/variables includes

Class	Strength - Peak (W/m ²)	What can they do to Earth?
B	$I < 10^{-6}$	Too small to harm Earth.
C	$10^{-6} \leq I < 10^{-5}$	Small with few noticeable consequences on Earth.
M	$10^{-5} \leq I < 10^{-4}$	Can cause brief radio blackouts that affect Earth's polar regions and minor radiation storms.
X	$I \geq 10^{-4}$	Can trigger planet-wide radio blackouts and long-lasting radiation storms.

Table 1: Classification of solar flares and their effects on Earth. This table is sourced from [Stanford Solar Center](#).

- Observation sequence $\mathbf{y}^{(l)} = (y_1^{(l)}, \dots, y_{T_l}^{(l)}) \in \mathbb{R}^{T_l}$, where T_l denotes the length.
- (Unobservable) latent state sequence $\mathbf{z}^{(l)} = (z_1^{(l)}, \dots, z_{T_l}^{(l)}) \in \mathbb{R}^{T_l}$.
- State transition matrix $\mathbf{A}^{(l)} \in \mathbb{R}^{K \times K}$, with $\mathbf{A}^{(l)}(i, j) = P(z_{t+1}^{(l)} = j \mid z_t^{(l)} = i)$, $1 \leq i, j \leq K$, $1 \leq t \leq T_l - 1$.

We assume the j -th rows of all $\mathbf{A}^{(l)}$'s are i.i.d. drawn from $\text{Dirichlet}(\boldsymbol{\lambda}_j)$, $1 \leq j \leq K$. The global parameters are

- $\boldsymbol{\lambda}_j \in \mathbb{R}^K$, $1 \leq j \leq K$.
- Observation matrix $\mathbf{B} \in \mathbb{R}^{K \times M}$, with $\mathbf{B}(j, k) = P(y_t^{(l)} = k \mid z_t^{(l)} = j)$, $1 \leq j \leq K$, $1 \leq k \leq M$, $1 \leq l \leq N$, $1 \leq t \leq T_l$.
- Initial state distribution $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K) \in \mathbb{R}^K$, with $\pi_k = P(z_1^{(l)} = k)$, $1 \leq l \leq N$, $1 \leq k \leq K$.

Assuming all trajectories are independent, the complete likelihood is give by

$$p\left(\left\{\mathbf{y}^{(n)}, \mathbf{z}^{(n)}\right\}_{n=1}^N \mid \left\{\mathbf{A}^{(n)}\right\}_{n=1}^N, \mathbf{B}, \boldsymbol{\pi}\right) = \prod_{l=1}^N \pi_{z_1^{(l)}} \prod_{t=2}^{T_l} \mathbf{A}^{(l)}(z_{t-1}^{(l)}, z_t^{(l)}) \prod_{t=1}^{T_l} \mathbf{B}(z_t^{(l)}, y_t^{(l)}). \quad (1)$$

For parameter estimation, we use Markov Chain Monte Carlo (MCMC; [Liu \(2001\)](#)) algorithms. Here are the prior distributions we adopt:

- For $\boldsymbol{\lambda}_j \in \mathbb{R}^K$ ($1 \leq j \leq K$): Independent components, each following $\text{Exp}(1)$.
- For \mathbf{B} : A flat prior with each row of \mathbf{B} drawn from $\text{Dirichlet}(\mathbf{1}_M)$.
- For $\boldsymbol{\pi}$: A flat prior $\text{Dirichlet}(\mathbf{1}_K)$.

We use a (group) Gibbs sampler to draw samples from the posterior distribution. The sampling process at each iteration can be stated as follows:

- For all $1 \leq j \leq K$, sample \mathbf{b}_j (the j -th row of \mathbf{B}) according to

$$p\left(\mathbf{b}_j \mid \left\{\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)}\right\}_{n=1}^N, \mathbf{B} \setminus \mathbf{b}_j, \boldsymbol{\pi}, \{\boldsymbol{\lambda}_k\}_{k=1}^K\right) \quad (2)$$

$$\propto p(\mathbf{b}_j) p\left(\left\{\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)}\right\}_{n=1}^N \mid \mathbf{B}, \boldsymbol{\pi}, \{\boldsymbol{\lambda}_k\}_{k=1}^K\right) \quad (3)$$

$$= p(\mathbf{b}_j) \prod_{n=1}^N p\left(\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)} \mid \mathbf{B}, \boldsymbol{\pi}, \{\boldsymbol{\lambda}_k\}_{k=1}^K\right) \quad (4)$$

$$\propto p(\mathbf{b}_j) \prod_{n=1}^N p\left(\mathbf{y}^{(n)}, \mathbf{z}^{(n)} \mid \mathbf{A}^{(n)}, \mathbf{B}, \boldsymbol{\pi}\right) \quad (5)$$

$$\propto p(\mathbf{b}_j) \prod_{n=1}^N \prod_{t=1}^{T_n} \left[b_j\left(y_t^{(n)}\right)\right]^{I\left(z_t^{(n)}=j\right)} \quad (6)$$

$$= \text{Dirichlet}\left(1 + \sum_{n=1}^N \sum_{t=1}^{T_n} I\left(z_t^{(n)} = j, y_t^{(n)} = 1\right), \dots, 1 + \sum_{n=1}^N \sum_{t=1}^{T_n} I\left(z_t^{(n)} = j, y_t^{(n)} = M\right)\right). \quad (7)$$

- Sample $\boldsymbol{\pi}$ from $\text{Dirichlet}\left(1 + \sum_{n=1}^N I\left(z_1^{(n)} = 1\right), \dots, 1 + \sum_{n=1}^N I\left(z_1^{(n)} = K\right)\right)$.

- For all $1 \leq l \leq N$ and $1 \leq j \leq K$, sample $\mathbf{a}_j^{(l)}$ (the j -th row of $\mathbf{A}^{(l)}$) from

$$\text{Dirichlet}\left(\lambda_j(1) + \sum_{t=2}^{T_l} I\left(z_{t-1}^{(l)} = j, z_t^{(l)} = 1\right), \dots, \lambda_j(K) + \sum_{t=2}^{T_l} I\left(z_{t-1}^{(l)} = j, z_t^{(l)} = K\right)\right)$$

- For all $1 \leq l \leq N$, sample $z_t^{(l)}$ sequentially according to

$$p\left(z_t^{(l)} = k \mid \mathbf{A}^{(l)}, \mathbf{B}, \boldsymbol{\pi}, \mathbf{y}^{(l)}, \mathbf{z}^{(l)} \setminus z_t^{(l)}\right) \propto a_{z_{t-1}^{(l)}}^{(l)}(k) a_k^{(l)}\left(z_{t+1}^{(l)}\right) b_k^{(l)}\left(y_t^{(l)}\right), \quad (8)$$

with appropriate adjustments for end effects:

$$p\left(z_1^{(l)} = k \mid \mathbf{A}^{(l)}, \mathbf{B}, \boldsymbol{\pi}, \mathbf{y}^{(l)}, \mathbf{z}^{(l)} \setminus z_1^{(l)}\right) \propto \pi_k a_k^{(l)}\left(z_2^{(l)}\right) b_k^{(l)}\left(y_1^{(l)}\right), \quad (9)$$

$$p\left(z_{T_l}^{(l)} = k \mid \mathbf{A}^{(l)}, \mathbf{B}, \boldsymbol{\pi}, \mathbf{y}^{(l)}, \mathbf{z}^{(l)} \setminus z_{T_l}^{(l)}\right) \propto a_{z_{T_l-1}^{(l)}}^{(l)}(k) b_k^{(l)}\left(y_{T_l}^{(l)}\right). \quad (10)$$

- For all $1 \leq l \leq K$, sample λ_l according to

$$p\left(\lambda_l \mid \left\{\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)}\right\}_{n=1}^N, \mathbf{B}, \boldsymbol{\pi}, \{\lambda_k\}_{k \neq l}\right) \quad (11)$$

$$\propto p(\lambda_l) p\left(\left\{\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)}\right\}_{n=1}^N \mid \mathbf{B}, \boldsymbol{\pi}, \{\lambda_k\}_{k=1}^K\right) \quad (12)$$

$$= p(\lambda_l) \prod_{n=1}^N p\left(\mathbf{A}^{(n)}, \mathbf{y}^{(n)}, \mathbf{z}^{(n)} \mid \mathbf{B}, \boldsymbol{\pi}, \{\lambda_k\}_{k=1}^K\right) \quad (13)$$

$$\propto p(\lambda_l) \prod_{n=1}^N p\left(\mathbf{A}^{(n)} \mid \{\lambda_k\}_{k=1}^K\right) \quad (14)$$

$$\propto p(\lambda_l) \prod_{n=1}^N p\left(\mathbf{a}_l^{(n)} \mid \lambda_l\right) \quad (15)$$

$$= \exp\left(-\sum_{j=1}^K \lambda_l(j)\right) \prod_{n=1}^N \left(\frac{\Gamma\left(\sum_{j=1}^K \lambda_l(j)\right)}{\prod_{j=1}^K \Gamma(\lambda_l(j))} \prod_{j=1}^K \left(a_l^{(n)}(j)\right)^{\lambda_l(j)-1}\right), \quad (16)$$

which is performed via the Metropolis-Hasting algorithm.

Note that the initializations for $\mathbf{A}^{(l)}$'s and $\mathbf{z}^{(l)}$'s are obtained via the Baum-Welch algorithm (Baum and Petrie (1966); Baum et al. (1970)).

For $K = 2$, the posterior mean of \mathbf{B} is

$$\begin{pmatrix} 0.9592601 & 0.04063102 & 0.0001088795 \\ 0.2440964 & 0.63732690 & 0.1185766887 \end{pmatrix} \left| \begin{array}{l} \text{Safe state} \\ \text{Dangerous state} \end{array} \right. . \quad (17)$$

The posterior mean of $\boldsymbol{\pi}$ is (0.2523393, 0.7476607). The posterior mean of λ_j 's is¹

$$\begin{pmatrix} 4.609977 & 0.510623 \\ 1.095575 & 1.047507 \end{pmatrix} \left| \begin{array}{l} \lambda_1 \\ \lambda_2 \end{array} \right. . \quad (18)$$

For $K = 3$, the posterior mean of \mathbf{B} is

$$\begin{pmatrix} 0.9707082 & 0.02898114 & 0.000310616 \\ 0.1139200 & 0.66911719 & 0.216962810 \\ 0.6706834 & 0.32390300 & 0.005413581 \end{pmatrix} \left| \begin{array}{l} \text{Safe state} \\ \text{Dangerous state} \\ \text{Moderate state} \end{array} \right. . \quad (19)$$

The posterior mean of $\boldsymbol{\pi}$ is (0.319113246, 0.671730366, 0.009156388). The posterior mean of λ_j 's is²

$$\begin{pmatrix} 7.6318634 & 0.03586825 & 0.6334422 \\ 0.4577618 & 0.34338912 & 0.4907412 \\ 0.8633875 & 1.12483292 & 5.1629052 \end{pmatrix} \left| \begin{array}{l} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{array} \right. . \quad (20)$$

Our results reveal the existence of a latent state which is more likely to produce strong flare events. In addition, this state is “active” in the sense that it tends to transition to other states, which is consistent with the physical intuition that high-energy states are usually unstable. Moreover, a flaring trajectory is more likely to begin in the high-energy state.

¹For example, the first rows of all $\mathbf{A}^{(l)}$'s can be considered as drawn from Dirichlet(4.609977, 0.510623).

²For example, the first rows of all $\mathbf{A}^{(l)}$'s can be considered as drawn from Dirichlet(7.6318634, 0.03586825, 0.6334422).

1.3 Lessons from the Internship

This research experience has not only deepened my understanding of multivariate methods but also broadened my perspective on real-world problem-solving. Specifically, I have gained the following lessons:

Real Data. Distribution shifts, misclassified instances, missing values, and weak signals (low signal-to-noise ratios) are commonly encountered in real-world data, which is totally different from “textbook examples”.

Attitudes toward Existing Literature. Everyone makes mistakes. Good researchers should not completely trust the results from a paper without careful thought. It is sometimes necessary to verify the claims through hands-on implementation.

Simulation Studies. When faced with challenges or counterintuitive results in real-data analysis, conducting simulation studies with synthetic data can often provide valuable insights and validate the reasoning.

Communication and Cooperation. Cooperation is essential in modern research. Good researchers should always save their collaborators’ time and communicate effectively. For example, preparing a brief document summarizing key points from the previous meeting, recent progress, and current questions can greatly facilitate the discussion—this is a habit I have maintained.

Academic Writing. In formal writing, it is important to avoid vague statements such as “almost”, “many (much)”, “most”, “roughly”, and “approximately”. Instead, we can use precise numbers to convey clear and accurate information. Additionally, taking high-quality notes can help clarify the challenges we are facing and the tools at our disposal. Summarizing the ideas in organized blocks (which is a habit I have maintained) can also save time and improve efficiency in academic writing.

2 Statistical Network Analysis

Preliminaries. We denote a network with n nodes by its adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, where $A_{ij} = 1$ if there is an edge (link) between nodes i and j and $A_{ij} = 0$ otherwise. We do not allow for self-edges (self-loops), i.e., the diagonal entries of \mathbf{A} are all zero.

2.1 The Stochastic Block Model

Among various models, the stochastic block model (SBM; [Holland et al. \(1983\)](#)) has attracted much attention and is arguably the best studied and most commonly used. Below, we briefly review its setting.

Suppose there are K communities. Each node belongs to only one of the communities. Let $\mathbf{c} = (c_1, \dots, c_n) \in \{1, \dots, K\}^n$ denote the true community labels of the nodes, and assume c_i ’s are i.i.d. categorical variables with parameter vector $\boldsymbol{\pi} = (\pi_1, \dots, \pi_K)$, where $\pi_k \geq 0$ for all $1 \leq k \leq K$ and $\sum_{k=1}^K \pi_k = 1$. Conditional on the community labels, the edge variables A_{ij} ’s are independent Bernoulli variables with $\mathbb{E}[A_{ij} \mid \mathbf{c}] = P(A_{ij} = 1 \mid \mathbf{c}) = P_{c_i c_j}$, where $\mathbf{P} \in [0, 1]^{K \times K}$

is the symmetric edge-probability matrix with the kl -entry P_{kl} characterizing the probability of connection between nodes in communities k and l . Let $\mathbf{\Omega} = (\boldsymbol{\pi}, \mathbf{P})$.

2.1.1 Why Estimation via EM is Intractable

It would be natural and logical to treat \mathbf{c} as latent variables and attempt to fit the SBM using EM algorithm. The log-likelihood for complete data (i.e., \mathbf{A} and \mathbf{c}) is computed as

$$\ell(\mathbf{\Omega}; \mathbf{A}, \mathbf{c}) = \log P(\mathbf{A}, \mathbf{c} \mid \mathbf{\Omega}) \quad (21)$$

$$= \log P(\mathbf{c} \mid \boldsymbol{\pi}) + \log P(\mathbf{A} \mid \mathbf{c}, \mathbf{P}) \quad (22)$$

$$= \sum_{i=1}^n \log P(c_i \mid \boldsymbol{\pi}) + \sum_{i < j} \log P(A_{ij} \mid c_i, c_j, \mathbf{P}) \quad (23)$$

$$= \sum_{i=1}^n \sum_{k=1}^K I(c_i = k) \log \pi_k \quad (24)$$

$$+ \sum_{i < j} [A_{ij} \log P(A_{ij} = 1 \mid c_i, c_j, \mathbf{P}) + (1 - A_{ij}) \log P(A_{ij} = 0 \mid c_i, c_j, \mathbf{P})] \quad (25)$$

$$= \sum_{i=1}^n \sum_{k=1}^K I(c_i = k) \log \pi_k \quad (26)$$

$$+ \sum_{i < j} \sum_{k=1}^K \sum_{l=1}^K I(c_i = k) I(c_j = l) [A_{ij} \log P_{kl} + (1 - A_{ij}) \log (1 - P_{kl})]. \quad (27)$$

In order to perform EM algorithm, we first need to derive the posterior distribution of \mathbf{c} , given the observed network \mathbf{A} and the model parameters $\mathbf{\Omega}$. Specifically, we have

$$P(\mathbf{c} \mid \mathbf{A}, \mathbf{\Omega}) = \frac{P(\mathbf{c}, \mathbf{A} \mid \mathbf{\Omega})}{P(\mathbf{A} \mid \mathbf{\Omega})} \quad (28)$$

$$\propto P(\mathbf{c}, \mathbf{A} \mid \mathbf{\Omega}) \quad (29)$$

$$= P(\mathbf{c} \mid \boldsymbol{\pi}) P(\mathbf{A} \mid \mathbf{c}, \mathbf{P}) \quad (30)$$

$$= \prod_{i=1}^n \pi_{c_i} \prod_{i < j} (P_{c_i c_j})^{A_{ij}} (1 - P_{c_i c_j})^{1 - A_{ij}}, \quad (31)$$

which leads to the following updating formula:

$$\pi_k^{(t+1)} \propto \sum_{i=1}^n P(c_i = k \mid \mathbf{A}, \mathbf{\Omega}^{(t)}), \quad 1 \leq k \leq K, \quad (32)$$

$$P_{kl}^{(t+1)} \propto \sum_{i < j} P(c_i = k, c_j = l \mid \mathbf{A}, \mathbf{\Omega}^{(t)}) A_{ij}, \quad 1 \leq k, l \leq K. \quad (33)$$

We note, however, that the posterior distribution of \mathbf{c} given by Eq. (31) is intractable for large n and K , since it is of order $\mathcal{O}(K^n)$. It is also computationally infeasible to solve for $\arg \max_{\mathbf{c}} P(\mathbf{c} \mid \mathbf{A}, \mathbf{\Omega})$, which poses a great challenge for community detection.

2.2 A Pseudo-Likelihood Method

Wang et al. (2023) proposed a pseudo-likelihood-based method³ which performs well for both small and large scale networks. Here we describe their method for fitting the SBM.

We introduce an initial column labeling vector $\mathbf{e} = (e_1, \dots, e_n) \in \{1, \dots, K\}^n$, and treat \mathbf{c} as a latent row labeling vector⁴. Let $\mathbf{a}_i = (A_{i1}, \dots, A_{in})$ denote the i -th row of \mathbf{A} , $1 \leq i \leq n$. The pseudo-likelihood function⁵ is defined as

$$\mathcal{L}_{\text{PL}}(\boldsymbol{\Omega}, \mathbf{e}; \{\mathbf{a}_i\}) = P(\{\mathbf{a}_i\} \mid \boldsymbol{\Omega}, \mathbf{e}) \quad (\text{Treat } \mathbf{e} \text{ as model parameters}) \quad (34)$$

$$= \prod_{i=1}^n P(\mathbf{a}_i \mid \boldsymbol{\Omega}, \mathbf{e}) \quad (\text{Ignore the potential dependency}) \quad (35)$$

$$= \prod_{i=1}^n \sum_{k=1}^K P(\mathbf{a}_i, c_i = k \mid \boldsymbol{\Omega}, \mathbf{e}) \quad (\text{Treat } \mathbf{c} \text{ as a vector of latent variables}) \quad (36)$$

$$= \prod_{i=1}^n \left\{ \sum_{k=1}^K P(c_i = k \mid \boldsymbol{\Omega}) P(\mathbf{a}_i \mid c_i = k, \boldsymbol{\Omega}, \mathbf{e}) \right\} \quad (37)$$

$$= \prod_{i=1}^n \left\{ \sum_{k=1}^K \pi_k \prod_{j=1}^n P(A_{ij} \mid c_i = k, \boldsymbol{\Omega}, \mathbf{e}) \right\} \quad (38)$$

$$= \prod_{i=1}^n \left\{ \sum_{k=1}^K \pi_k \prod_{j=1}^n P_{ke_j}^{A_{ij}} (1 - P_{ke_j})^{1-A_{ij}} \right\}, \quad (39)$$

with its logarithm as

$$\ell_{\text{PL}}(\boldsymbol{\Omega}, \mathbf{e}; \{\mathbf{a}_i\}) = \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \pi_k \prod_{j=1}^n P_{ke_j}^{A_{ij}} (1 - P_{ke_j})^{1-A_{ij}} \right\}. \quad (40)$$

2.2.1 Estimation via Alternating Update

We use an alternating update algorithm to estimate the model parameters $\boldsymbol{\Omega}$ and \mathbf{e} .

Given the current estimate of column labeling vector denoted by $\hat{\mathbf{e}}$, maximizing $\ell_{\text{PL}}(\boldsymbol{\Omega}, \hat{\mathbf{e}}; \{\mathbf{a}_i\})$ over $\boldsymbol{\Omega}$ can be conducted by a vanilla EM algorithm (treating c_i 's as latent variables), where each M-step admits a closed form:

$$\pi_k^{(t+1)} \propto \sum_{i=1}^n P(c_i = k \mid \mathbf{a}_i, \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}}), \quad 1 \leq k \leq K, \quad (41)$$

$$P_{kl}^{(t+1)} \propto \sum_{i=1}^n \sum_{j=1}^n A_{ij} P(c_i = k \mid \mathbf{a}_i, \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}}) I(\hat{e}_j = l), \quad 1 \leq k, l \leq K. \quad (42)$$

³The idea of pseudo-likelihood dates back to Besag (1974) and generally involves ignoring some of the dependency structure in the data to simplify the likelihood and make it more tractable.

⁴Here we are lifting the symmetry constraint on the adjacency matrix \mathbf{A} .

⁵In the original paper Wang et al. (2023), this function is referred to as the *profile-pseudo likelihood*, as $\boldsymbol{\Omega}$ is a nuisance parameter while \mathbf{e} is the parameter of interest.

The quantity $P(c_i = k \mid \mathbf{a}_i, \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}})$ can be computed as

$$P(c_i = k \mid \mathbf{a}_i, \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}}) \propto P(c_i = k, \mathbf{a}_i \mid \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}}) \quad (43)$$

$$= P(c_i = k \mid \boldsymbol{\Omega}^{(t)}) P(\mathbf{a}_i \mid c_i = k, \boldsymbol{\Omega}^{(t)}, \hat{\mathbf{e}}) \quad (44)$$

$$= \pi_k^{(t)} \prod_{j=1}^n \left(P_{k\hat{e}_j}^{(t)} \right)^{A_{ij}} \left(1 - P_{k\hat{e}_j}^{(t)} \right)^{1-A_{ij}}. \quad (45)$$

Given the current estimate of $\boldsymbol{\Omega}$ denoted by $\hat{\boldsymbol{\Omega}}$, maximizing $\ell_{\text{PL}}(\hat{\boldsymbol{\Omega}}, \mathbf{e}; \{\mathbf{a}_i\})$ over all possible \mathbf{e} 's is of order $\mathcal{O}(K^n)$, which is in fact intractable. In order to efficiently update \mathbf{e} , Wang et al. (2023) proposed the following updating formula:

$$e_j^{(s+1)} \leftarrow \arg \max_{k \in \{1, \dots, K\}} \sum_{i=1}^n \sum_{l=1}^K P(c_i = l \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}) \left\{ A_{ij} \log \hat{P}_{lk} + (1 - A_{ij}) \log(1 - \hat{P}_{lk}) \right\}, \quad (46)$$

where the update for \mathbf{e} is obtained separately for each node. Eq. (46) can be intuitively justified as follows:

$$\sum_{i=1}^n \sum_{l=1}^K P(c_i = l \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}) \left\{ A_{ij} \log \hat{P}_{le_j^{(s+1)}} + (1 - A_{ij}) \log(1 - \hat{P}_{le_j^{(s+1)}}) \right\} \quad (47)$$

$$= \sum_{i=1}^n \sum_{l=1}^K P(c_i = l \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}) \log P(A_{ij} \mid c_i = l, e_j^{(s+1)}, \hat{\boldsymbol{\Omega}}) \quad (48)$$

$$= \sum_{i=1}^n \mathbb{E}_{c_i \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}} \left[\log P(A_{ij} \mid c_i, e_j^{(s+1)}, \hat{\boldsymbol{\Omega}}) \right] \quad (49)$$

$$\leq \sum_{i=1}^n \log \left\{ \mathbb{E}_{c_i \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}} \left[P(A_{ij} \mid c_i, e_j^{(s+1)}, \hat{\boldsymbol{\Omega}}) \right] \right\}. \quad (50)$$

We can also formally prove $\ell_{\text{PL}}(\hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s+1)}; \{\mathbf{a}_i\}) \geq \ell_{\text{PL}}(\hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}; \{\mathbf{a}_i\})$. To simplify the notation, we write

$$\hat{\tau}_{ik}^{(s)} = P(c_i = k \mid \mathbf{a}_i, \hat{\boldsymbol{\Omega}}, \mathbf{e}^{(s)}) \quad (51)$$

$$= \frac{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}} \right)^{1-A_{ij}}}{\sum_{k=1}^K \hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}} \right)^{1-A_{ij}}}. \quad (52)$$

Then we have

$$\ell_{\text{PL}}(\hat{\Omega}, \mathbf{e}^{(s+1)}; \{\mathbf{a}_i\}) - \ell_{\text{PL}}(\hat{\Omega}, \mathbf{e}^{(s)}; \{\mathbf{a}_i\}) \quad (53)$$

$$= \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}} \right\} - \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}} \right\} \quad (54)$$

$$= \sum_{i=1}^n \log \left\{ \frac{\sum_{k=1}^K \hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}}}{\sum_{k=1}^K \hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}} \right\} \quad (55)$$

$$= \sum_{i=1}^n \log \left\{ \frac{\sum_{k=1}^K \frac{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}}}{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}}}{\sum_{k=1}^K \frac{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}}{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}}} \right\} \quad (56)$$

$$= \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \frac{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}}}{\hat{\pi}_k \prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}} \hat{\tau}_{ik}^{(s)} \right\} \quad (57)$$

$$\geq \sum_{i=1}^n \sum_{k=1}^K \hat{\tau}_{ik}^{(s)} \log \left\{ \frac{\prod_{j=1}^n \hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}}}{\prod_{j=1}^n \hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}}} \right\} \quad (58)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \sum_{j=1}^n \hat{\tau}_{ik}^{(s)} \left\{ \log \left[\hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}} \right] - \log \left[\hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}} \right] \right\} \quad (59)$$

$$= \sum_{j=1}^n \left(\sum_{i=1}^n \sum_{k=1}^K \hat{\tau}_{ik}^{(s)} \left\{ \log \left[\hat{P}_{ke_j^{(s+1)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s+1)}}\right)^{1-A_{ij}} \right] - \log \left[\hat{P}_{ke_j^{(s)}}^{A_{ij}} \left(1 - \hat{P}_{ke_j^{(s)}}\right)^{1-A_{ij}} \right] \right\} \right) \quad (60)$$

$$\geq 0, \quad (61)$$

where the first inequality is due to Jensen's inequality, and the second inequality is due to the updating formula (46).

Although the alternating update algorithm stated above is not guaranteed to maximize the pseudo log-likelihood $\ell_{\text{PL}}(\Omega, \mathbf{e}; \{\mathbf{a}_i\})$, it ensures a non-negative increment in $\ell_{\text{PL}}(\Omega, \mathbf{e}; \{\mathbf{a}_i\})$ at each iteration.

When the algorithm converges, the output $\hat{\mathbf{e}}$ gives the result of community detection.

2.3 Jin's SCORE Algorithm

Jin (2015) proposed the Spectral Clustering On Ratios-of-Eigenvectors (SCORE) algorithm for fast community detection in the degree-corrected SBM (to be introduced in Section 2.4). Here

we discuss the intuition behind SCORE in the context of SBM⁶.

We treat $\mathbf{c} = (c_1, \dots, c_n) \in \{1, \dots, K\}^n$ as unknown model parameters, rather than as random variables. For each node i , we write the column vector $\mathbf{e}_i \in \mathbb{R}^K$ as its membership vector where $e_i(k) = I(c_i = k)$ for all $1 \leq k \leq K$. To write the model in matrix form, we define $\mathbf{E} := (\mathbf{e}_1, \dots, \mathbf{e}_n)^\top = [\mathbf{E}_1, \dots, \mathbf{E}_K] \in \mathbb{R}^{n \times K}$, and write $\mathbf{\Omega} := \mathbf{E} \mathbf{P} \mathbf{E}^\top = \sum_{i=1}^K \sum_{j=1}^K P_{ij} \mathbf{E}_i \mathbf{E}_j^\top \in \mathbb{R}^{n \times n}$. All the other notations introduced earlier will be retained. Our objective is to recover \mathbf{c} from the observed network \mathbf{A} .

In SBM, for all edge variables $\{A_{ij}, 1 \leq i < j \leq n\}$, we have

$$\mathbb{E}[A_{ij}] = P(A_{ij} = 1) \tag{62}$$

$$= P_{c_i c_j} \quad (\text{Treat } \mathbf{c} \text{ as model parameters}) \tag{63}$$

$$= \mathbf{e}_i^\top \mathbf{P} \mathbf{e}_j \quad (\text{Definition of } \mathbf{e}_i\text{'s}) \tag{64}$$

$$= (\mathbf{E} \mathbf{P} \mathbf{E}^\top)(i, j) \quad (\text{Definition of } \mathbf{E}) \tag{65}$$

$$= \mathbf{\Omega}(i, j). \quad (\text{Definition of } \mathbf{\Omega}) \tag{66}$$

Since we do not allow for self-edges, we have $\mathbb{E}[\mathbf{A}] = \mathbf{\Omega} - \text{diag}(\mathbf{\Omega})$. It follows that

$$\mathbf{A} = \mathbb{E}[\mathbf{A}] + (\mathbf{A} - \mathbb{E}[\mathbf{A}]) \tag{67}$$

$$= \mathbf{\Omega} - \text{diag}(\mathbf{\Omega}) + \mathbf{W}. \quad (\text{Define } \mathbf{W} = \mathbf{A} - \mathbb{E}[\mathbf{A}]) \tag{68}$$

Recall that our goal is to recover \mathbf{c} from \mathbf{A} . We ought to find out “which part of the data contains the information” (Tukey (1965)). In this context, we note the following:

- $\mathbf{\Omega}$ contains the most direct information of the community labels.
- $\text{diag}(\mathbf{\Omega})$ contains negligible information compared to $\mathbf{\Omega}$.
- \mathbf{W} should be considered as some noise, since the upper triangular part of \mathbf{W} consists of independent (but not identical) centered-Bernoulli variables.

We next consider the oracle case where we have access to the main signal matrix $\mathbf{\Omega} = \mathbf{E} \mathbf{P} \mathbf{E}^\top = \sum_{i=1}^K \sum_{j=1}^K P_{ij} \mathbf{E}_i \mathbf{E}_j^\top$.

2.3.1 The Oracle Case

Suppose $\boldsymbol{\eta} \in \mathbb{R}^n$ is an eigenvector of $\mathbf{\Omega}$ with eigenvalue λ . We wish to find the connection between $\boldsymbol{\eta}$ and the community label vector \mathbf{c} .

The analysis starts with

$$\lambda \boldsymbol{\eta} = \mathbf{\Omega} \boldsymbol{\eta} \tag{69}$$

$$= \sum_{i=1}^K \sum_{j=1}^K P_{ij} \mathbf{E}_i \mathbf{E}_j^\top \boldsymbol{\eta}. \tag{70}$$

⁶In fact, Jin (2015) demonstrated both theoretically and empirically that the effect of degree heterogeneity can be mitigated by SCORE. In this note, however, we only consider a simpler model, the SBM, to highlight the underlying principles.

Then, for all $1 \leq l \leq K$, we have

$$\lambda \mathbf{E}_l^\top \boldsymbol{\eta} = \mathbf{E}_l^\top \sum_{i=1}^K \sum_{j=1}^K P_{ij} \mathbf{E}_i \mathbf{E}_j^\top \boldsymbol{\eta} \quad (71)$$

$$= \|\mathbf{E}_l\|_2^2 \sum_{j=1}^K P_{lj} \mathbf{E}_j^\top \boldsymbol{\eta}. \quad (\mathbf{E} \text{ has orthogonal columns}) \quad (72)$$

We may rewrite the previous line in its matrix form:

$$\text{diag}(\|\mathbf{E}_1\|_2^2, \dots, \|\mathbf{E}_K\|_2^2) \mathbf{P} \begin{pmatrix} \mathbf{E}_1^\top \boldsymbol{\eta} \\ \vdots \\ \mathbf{E}_K^\top \boldsymbol{\eta} \end{pmatrix} = \begin{pmatrix} \|\mathbf{E}_1\|_2^2 \sum_{j=1}^K P_{1j} \mathbf{E}_j^\top \boldsymbol{\eta} \\ \vdots \\ \|\mathbf{E}_K\|_2^2 \sum_{j=1}^K P_{Kj} \mathbf{E}_j^\top \boldsymbol{\eta} \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{E}_1^\top \boldsymbol{\eta} \\ \vdots \\ \mathbf{E}_K^\top \boldsymbol{\eta} \end{pmatrix}. \quad (73)$$

To simplify the notation, we write $\mathbf{D} = \text{diag}(\|\mathbf{E}_1\|_2^2, \dots, \|\mathbf{E}_K\|_2^2)$. It follows that $(\mathbf{E}_1^\top \boldsymbol{\eta}, \dots, \mathbf{E}_K^\top \boldsymbol{\eta})^\top$ is an eigenvector of $\mathbf{D}\mathbf{P}$ with eigenvalue λ . Further, if $\mathbf{a} \in \mathbb{R}^K$ is an eigenvector of $\mathbf{D}\mathbf{P}$, then $\boldsymbol{\eta}_{\mathbf{a}}$ is an eigenvector of $\boldsymbol{\Omega}$ with the same eigenvalue, where

$$\boldsymbol{\eta}_{\mathbf{a}} = \sum_{l=1}^K \frac{a(l)}{\|\mathbf{E}_l\|_2} \mathbf{E}_l. \quad (74)$$

We assume that $\mathbf{D}\mathbf{P}$ has eigenvectors $\mathbf{a}_k, 1 \leq k \leq K$. Then the eigenvectors of $\boldsymbol{\Omega}$ are given by⁷

$$\boldsymbol{\eta}_k = \sum_{l=1}^K \frac{a_k(l)}{\|\mathbf{E}_l\|_2} \mathbf{E}_l, \quad 1 \leq k \leq K, \quad (75)$$

which implies that, for $1 \leq i \leq n$, we have

$$\eta_k(i) = \sum_{l=1}^K \frac{a_k(l)}{\|\mathbf{E}_l\|_2} E_l(i) \quad (76)$$

$$= \sum_{l=1}^K \frac{a_k(l)}{\|\mathbf{E}_l\|_2} I(c_i = l) \quad (77)$$

$$= \frac{a_k(c_i)}{\|\mathbf{E}_{c_i}\|_2}. \quad (78)$$

This suggests that $\eta_k(i) = \eta_k(j)$ for all $1 \leq k \leq K$ if $c_i = c_j$ (i.e., nodes i and j belong to the same community). We may take the ratio⁸,

$$R(i, k) = \frac{\eta_{k+1}(i)}{\eta_1(i)}, \quad 1 \leq i \leq n, 1 \leq k \leq K-1, \quad (79)$$

$$= \frac{a_{k+1}(c_i)}{a_1(c_i)}, \quad (80)$$

and define the matrix $\mathbf{R} \in \mathbb{R}^{n \times (K-1)}$ with its (i, k) entry as $R(i, k)$. We finally conclude that

⁷The matrix $\boldsymbol{\Omega}$ has at most K eigenvalues since $\text{rank}(\boldsymbol{\Omega}) = \text{rank}(\mathbf{E}\mathbf{P}\mathbf{E}^\top) \leq \min\{\text{rank}(\mathbf{E}), \text{rank}(\mathbf{P})\} \leq K$.

⁸When using the SBM, it is unnecessary to take the ratio. However, in the degree-corrected SBM, the matrix \mathbf{E} should be redefined as $\mathbf{E} = (\theta_1 \mathbf{e}_1, \dots, \theta_n \mathbf{e}_n)^\top$, where the ratios eliminate the nuisance degree parameters (i.e., θ_i 's).

- The matrix \mathbf{R} has only K distinct rows.
- For any $1 \leq i \neq j \leq n$ with $c_i = c_j$ (i.e., nodes i and j belong to the same community), the i -th and j -th rows of \mathbf{R} are identical.

We are now ready to present the SCORE community detection algorithm in the oracle case:

- Compute the eigenvectors of $\mathbf{\Omega}$ and denote them by $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_K$.
- Define the matrix $\mathbf{R} \in \mathbb{R}^{n \times (K-1)}$ with its (i, k) entry as $R(i, k) = \eta_{k+1}(i)/\eta_1(i)$.
- Nodes i and j belong to the same community if the i -th and j -th rows of \mathbf{R} are identical.

2.3.2 The Real Case

We now consider the real case where \mathbf{A} , instead of $\mathbf{\Omega}$, is observed.

Considering what we have learned from the oracle case, it is natural and logical to write the algorithm as follows:

- Compute the K (unit-norm) leading eigenvectors⁹ of \mathbf{A} and denote them by $\hat{\boldsymbol{\eta}}_1, \dots, \hat{\boldsymbol{\eta}}_K$.
- Define the matrix $\hat{\mathbf{R}} \in \mathbb{R}^{n \times (K-1)}$ with its (i, k) entry as $\hat{R}(i, k) = \hat{\eta}_{k+1}(i)/\hat{\eta}_1(i)$.
- Apply a clustering algorithm to $\hat{\mathbf{R}}$, such as k -means¹⁰.

2.4 Extensions of the SBM

We introduce several extensions of the SBM. In the following paragraphs, all notations introduced previously will be retained and used consistently unless explicitly stated otherwise or there exists a conflict. Throughout, we assume that there are K communities in total.

The mixed-membership SBM model (MMSBM; [Airoldi et al. \(2008\)](#)) allows each node to belong to multiple communities by introducing a weight vector $\boldsymbol{\pi}_i \in \mathbb{R}^K$ for each node i . Conditional on $\{\boldsymbol{\pi}_i, 1 \leq i \leq n\}$, the edge variables $\{A_{ij}, 1 \leq i < j \leq n\}$ are assumed to be independent Bernoulli variables with

$$\mathbb{E}[A_{ij} \mid \boldsymbol{\pi}_i, \boldsymbol{\pi}_j] = P(A_{ij} = 1 \mid \boldsymbol{\pi}_i, \boldsymbol{\pi}_j) \quad (81)$$

$$= \sum_{k=1}^K \sum_{l=1}^K \pi_i(k) P_{kl} \pi_j(l) \quad (82)$$

$$= \boldsymbol{\pi}_i^\top \mathbf{P} \boldsymbol{\pi}_j. \quad (83)$$

It is noteworthy that in the original paper [Airoldi et al. \(2008\)](#), the weight vectors (i.e., $\boldsymbol{\pi}_i$'s) are treated as independent random variables drawn from a Dirichlet distribution denoted by $\text{Dir}(\boldsymbol{\alpha})$. The parameters to be estimated are \mathbf{P} and $\boldsymbol{\alpha}$.

⁹By “leading eigenvectors”, we are comparing the absolute values of the eigenvalues.

¹⁰The k -means algorithm is also the choice in the original paper [Jin \(2015\)](#).

The degree-corrected SBM (DCSBM; [Karrer and Newman \(2011\)](#)) accounts for degree heterogeneity by incorporating additional degree parameters. Specifically, conditional on the label vector \mathbf{c} , DCSBM assumes that the edge variables $\{A_{ij}, 1 \leq i < j \leq n\}$ are independent Poisson variables with

$$\mathbb{E}[A_{ij} \mid \mathbf{c}] = \theta_i \theta_j \lambda_{c_i c_j}, \quad (84)$$

where $\mathbf{\Lambda} = [\lambda_{kl}]$ is a $K \times K$ symmetric matrix and $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ is a degree parameter vector.

The Degree-Corrected Mixed-Membership (DCMM) model ([Jin et al. \(2017\)](#); [Ji et al. \(2022\)](#); [Jin et al. \(2024\)](#)) permits both degree heterogeneity and mixed memberships. Specifically, DCMM assumes that each node i is associated with a K -dimensional weight vector $\boldsymbol{\pi}_i \in \mathbb{R}^K$ where for $1 \leq k \leq K$,

$$\pi_i(k) = \text{the } k\text{-th component of } \boldsymbol{\pi}_i \quad (85)$$

$$= \text{the fractional weight of node } i \text{ on community } k. \quad (86)$$

Below is an example from [Ji et al. \(2022\)](#) which helps clarify the meaning of $\boldsymbol{\pi}_i$ in DCMM:

“Suppose $K = 3$ and we have three communities, each being a primary area in statistics: ‘Bayes’, ‘Biostatistics’, and ‘Non-parametric’. Suppose for author i , $\boldsymbol{\pi}_i = (0.5, 0.3, 0.2)^\top$. In this case, we think author i has 50%, 30%, and 20% of his research interest or impact in these primary areas, respectively.”

DCMM further assumes that the edge variables $\{A_{ij}, 1 \leq i < j \leq n\}$ are independent Bernoulli variables with

$$\mathbb{E}[A_{ij}] = P(A_{ij} = 1) \quad (87)$$

$$= \theta_i \theta_j \sum_{k=1}^K \sum_{l=1}^K \pi_i(k) P_{kl} \pi_j(l) \quad (88)$$

$$= \theta_i \theta_j \boldsymbol{\pi}_i^\top \mathbf{P} \boldsymbol{\pi}_j. \quad (89)$$

Note that in [Jin et al. \(2017\)](#), [Ji et al. \(2022\)](#), and [Jin et al. \(2024\)](#), $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$, $\{\boldsymbol{\pi}_i, 1 \leq i \leq n\}$, and \mathbf{P} are all treated as unknown model parameters to be estimated, rather than as random variables drawn from some distributions.

2.5 Prediction Models for Network-Linked Data

Apart from fitting a realistic and mathematically tractable model for an observed network, it is also interesting to predict a response variable from covariates by utilizing the network information. Here, we briefly introduce the work of [Li et al. \(2019\)](#), which adds a penalty on individual node effects to encourage similarity between predictions for linked nodes.

Set-up and Notation. The data $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$ consist of n observations, where $y_i \in \mathbb{R}$ is the response variable and $\mathbf{x}_i \in \mathbb{R}^p$ is the vector of covariates for observation i . We write $\mathbf{Y} = (y_1, \dots, y_n)^\top$ for the n -dimensional response vector, and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$ for the $n \times p$

design matrix. The adjacency matrix \mathbf{A} is defined as in previous sections. We write $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ for the *degree matrix*, where $d_i = \sum_{j=1}^n A_{ij}$ is the number of nodes connected to node i for $1 \leq i \leq n$. The (unnormalized) Laplacian is then given by $\mathbf{L} = \mathbf{D} - \mathbf{A}$.

2.5.1 Linear Regression with Network Cohesion

The simplest prediction model is perhaps the linear regression model¹¹:

$$\mathbf{Y} = \boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (90)$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top \in \mathbb{R}^n$ is the vector of individual node effects, and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^p$ is the vector of regression coefficients. Li et al. (2019) proposed the following loss function:

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{L} \boldsymbol{\alpha}, \quad (91)$$

where $\|\cdot\|_2$ is the L_2 norm and $\lambda > 0$ is a tuning parameter. The penalty term $\lambda \boldsymbol{\alpha}^\top \mathbf{L} \boldsymbol{\alpha}$ can be understood via the following derivation:

$$\boldsymbol{\alpha}^\top \mathbf{L} \boldsymbol{\alpha} = \boldsymbol{\alpha}^\top \mathbf{D} \boldsymbol{\alpha} - \boldsymbol{\alpha}^\top \mathbf{A} \boldsymbol{\alpha} \quad (\mathbf{L} = \mathbf{D} - \mathbf{A}) \quad (92)$$

$$= \sum_{i=1}^n \alpha_i^2 d_i - \sum_{i=1}^n \sum_{j=1}^n \alpha_i A_{ij} \alpha_j \quad (93)$$

$$= \sum_{i=1}^n \alpha_i^2 \left(\sum_{j=1}^n A_{ij} \right) - \sum_{i=1}^n \sum_{j=1}^n \alpha_i A_{ij} \alpha_j \quad (94)$$

$$= \sum_{i=1}^n \sum_{j=1}^n \alpha_i^2 A_{ij} - \sum_{i=1}^n \sum_{j=1}^n \alpha_i A_{ij} \alpha_j \quad (95)$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i^2 + \alpha_j^2) A_{ij} - \sum_{i=1}^n \sum_{j=1}^n \alpha_i A_{ij} \alpha_j \quad (\mathbf{A} = \mathbf{A}^\top) \quad (96)$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i^2 + \alpha_j^2 - 2\alpha_i \alpha_j) A_{ij} \quad (97)$$

$$= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \alpha_j)^2 A_{ij}. \quad (98)$$

Thus, Eq. (91) penalizes differences between individual effects of nodes connected by an edge in the network.

Eq. (91) can be written in the matrix form:

$$L(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\alpha}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \mathbf{L} \boldsymbol{\alpha} \quad (99)$$

$$= \left\| \mathbf{Y} - (\mathbf{I}_n, \mathbf{X}) \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix} \right\|_2^2 + \lambda \begin{pmatrix} \boldsymbol{\alpha}^\top & \boldsymbol{\beta}^\top \end{pmatrix} \begin{pmatrix} \mathbf{L} & \mathbf{0}_{n \times p} \\ \mathbf{0}_{p \times n} & \mathbf{0}_{p \times p} \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix}, \quad (100)$$

¹¹In the original paper Li et al. (2019), generalized linear models and Cox's proportional hazard model (Cox (1972)) are also considered. However, the theoretic guarantee is provided only for linear models. In this note, we focus on the linear regression model to highlight the intuition.

which can be viewed as a ridge regression. To simplify the notation, we write

$$\tilde{\mathbf{X}} = (\mathbf{I}_n, \mathbf{X}), \quad (101)$$

$$\boldsymbol{\theta} = \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix}, \quad (102)$$

$$\mathbf{M} = \begin{pmatrix} \mathbf{L} & \mathbf{0}_{n \times p} \\ \mathbf{0}_{p \times n} & \mathbf{0}_{p \times p} \end{pmatrix}. \quad (103)$$

The minimizer of Eq. (91) admits the following closed form:

$$\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\alpha}}^\top, \hat{\boldsymbol{\beta}}^\top)^\top \quad (104)$$

$$= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{M})^{-1} \tilde{\mathbf{X}}^\top \mathbf{Y}, \quad (105)$$

provided that the matrix $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{M}$ is invertible.

In fact, we have

$$\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{M} = \begin{pmatrix} \mathbf{I}_n + \lambda \mathbf{L} & \mathbf{X} \\ \mathbf{X}^\top & \mathbf{X}^\top \mathbf{X} \end{pmatrix}. \quad (106)$$

Thus, if we assume that $\mathbf{X}^\top \mathbf{X}$ is invertible, then the matrix $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} + \lambda \mathbf{M}$ is invertible if and only if the **Schur complement** $\mathbf{I}_n + \lambda \mathbf{L} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \lambda \mathbf{L} + \mathbf{P}_{\mathbf{X}^\perp}$ is invertible. In practice, we can replace \mathbf{L} with $\mathbf{L} + \gamma \mathbf{I}_n$ to ensure numerical stability, where γ is a small positive constant.

2.5.2 A Bayesian Interpretation

The loss function (91) has a Bayesian interpretation. Specifically, we consider the following model:

$$\mathbf{Y} \mid \mathbf{X}, \boldsymbol{\alpha}, \boldsymbol{\beta} \sim \mathcal{N}(\boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n), \quad \boldsymbol{\alpha} \sim \pi_{\boldsymbol{\alpha}}(\boldsymbol{\Phi}), \quad \boldsymbol{\beta} \sim \pi_{\boldsymbol{\beta}}(\boldsymbol{\phi}), \quad (107)$$

where

- $\pi_{\boldsymbol{\alpha}}(\boldsymbol{\Phi})$ is the prior for $\boldsymbol{\alpha}$ with hyperparameter $\boldsymbol{\Phi}$,
- $\pi_{\boldsymbol{\beta}}(\boldsymbol{\phi})$ is the prior for $\boldsymbol{\beta}$ with hyperparameter $\boldsymbol{\phi}$,
- $\sigma^2 > 0$ is assumed to be known.

To reflect the lack of prior information about the regression coefficients, we specify a uniform prior on $\boldsymbol{\beta}$, i.e.,

$$\pi_{\boldsymbol{\beta}}(\boldsymbol{\phi}) \propto 1. \quad (108)$$

The posterior distribution of $\boldsymbol{\theta} = (\boldsymbol{\alpha}^\top, \boldsymbol{\beta}^\top)^\top$ is then given by

$$p(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{Y}) \propto p(\boldsymbol{\theta}) p(\mathbf{Y} \mid \mathbf{X}, \boldsymbol{\theta}) \quad (109)$$

$$\propto \pi_{\boldsymbol{\alpha}}(\boldsymbol{\Phi}) \pi_{\boldsymbol{\beta}}(\boldsymbol{\phi}) \exp \left(-\frac{1}{2\sigma^2} \left\| \mathbf{Y} - (\mathbf{I}_n, \mathbf{X}) \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix} \right\|_2^2 \right) \quad (110)$$

$$\propto \pi_{\boldsymbol{\alpha}}(\boldsymbol{\Phi}) \exp \left(-\frac{1}{2\sigma^2} \left\| \mathbf{Y} - \tilde{\mathbf{X}} \boldsymbol{\theta} \right\|_2^2 \right). \quad (111)$$

Thus, if we assume a Gaussian Markov random field prior on α with $\pi_\alpha(\Phi) = \mathcal{N}(\mathbf{0}_n, \Phi)$, where $\Phi = (\mathbf{L} + \gamma \mathbf{I}_n)^{-1}$, then the posterior mode is equivalent to the minimizer of Eq. (91).

2.5.3 Prediction and Choosing the Tuning Parameter

For in-sample prediction, we can simply use $\hat{\alpha} + \mathbf{X}\hat{\beta}$. In the original paper Li et al. (2019), the out-of-sample prediction task is to make predictions on a group of new subjects whose covariates as well as network connections (but not responses) become available after the model is fitted on training data.

Suppose there are n training samples and n' test samples. The Laplacian for the entire graph can be written as

$$\mathbf{L}' = \begin{pmatrix} \mathbf{L}_{11} & \mathbf{L}_{12} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{pmatrix}, \quad (112)$$

where \mathbf{L}_{11} corresponds to the n training samples and \mathbf{L}_{22} corresponds to the n' test samples. Similarly, we write the individual effect vector as $(\alpha_1^\top, \alpha_2^\top)^\top$, where $\alpha_1 = \hat{\alpha}$ is estimated from training data and α_2 needs to be predicted. Li et al. (2019) suggested we take advantage of cohesion and estimate α_2 by

$$\hat{\alpha}_2 = \arg \min_{\alpha_2} (\hat{\alpha}^\top, \alpha_2^\top) \mathbf{L}' (\hat{\alpha}^\top, \alpha_2^\top)^\top \quad (113)$$

$$= \arg \min_{\alpha_2} \left(\hat{\alpha}^\top \mathbf{L}_{12} \alpha_2 + \alpha_2^\top \mathbf{L}_{21} \hat{\alpha} + \alpha_2^\top \mathbf{L}_{22} \alpha_2 \right) \quad (114)$$

$$= \arg \min_{\alpha_2} \left(2\alpha_2^\top \mathbf{L}_{21} \hat{\alpha} + \alpha_2^\top \mathbf{L}_{22} \alpha_2 \right) \quad (115)$$

$$= -\mathbf{L}_{22}^{-1} \mathbf{L}_{21} \hat{\alpha}. \quad (116)$$

In addition, Li et al. (2019) reported that the naive cross-validation (random splitting) is appropriate for determining the tuning parameter $\lambda > 0$. This is likely because the problem is fundamentally a regression task, and we do not make inferences about the network structure¹².

3 Statistical Text Analysis

Ludwig Wittgenstein, the famous Austrian philosopher, once said,

“The limits of my language mean the limits of my world.”

Given n documents written with a vocabulary of p words, we denote the sets of all documents and words by $\{d^{(1)}, \dots, d^{(n)}\}$ and $\{w^{(1)}, \dots, w^{(p)}\}$, respectively. We assume that there are only K ($\ll \min\{n, p\}$) topics that are discussed by all these documents. We denote the set of all topics by $\{z^{(1)}, \dots, z^{(K)}\}$.

We write $X \in \mathbb{R}^{p \times n}$ as the *word-document-count* matrix, where $X(j, i)$ is the count of the j -th vocabulary in document i . We denote the word count vector of document i by $x_i = (X(1, i), \dots, X(p, i))^\top \in \mathbb{R}^p$. Then we have $X = [x_1, \dots, x_n]$.

¹²Recall that when making out-of-sample predictions, we assume access to the ground-truth network structure.

Throughout, we consider only *bag-of-words* models, i.e., models that focus on the counts of individual words, while neglecting word order and context.

3.1 Mixture of Unigrams

Nigam et al. (2000) used a topic model that Blei et al. (2003) referred to as the *mixture of unigrams* model¹³. We present the model here in a slightly modified form. For each topic $z^{(k)}$, we write its population word frequency vector as $A_k \in \mathbb{R}^p$. We assume that for all $1 \leq i \leq n$, document i is generated under some latent topic $z_i \in \{z^{(1)}, \dots, z^{(K)}\}$, and that

$$P(z_i = z^{(k)}) = w(k), \quad 1 \leq k \leq K. \quad (117)$$

Note that here $P(z_i = z^{(k)})$ does not vary with i , and each document exhibits exactly one topic¹⁴. Suppose document i has a total of N_i words. We further assume that

$$x_i \mid z_i = z^{(k)} \sim \text{Multinomial}(N_i, A_k), \quad 1 \leq i \leq n, \quad (118)$$

that is, once the topic $z_i = z^{(k)}$ is given, the distribution of x_i is multinomial with parameters (N_i, A_k) .

We write $w = (w(1), \dots, w(K))^\top \in \mathbb{R}^K$ to combine the weights of different topics into a vector, which we refer to as the *topic weight vector*. We also refer to $A = [A_1, \dots, A_K] \in \mathbb{R}^{p \times K}$ as the *topic matrix*. Then we have $X = [x_1, \dots, x_n]$. The main focus is to estimate $\theta = (A, w)$ from X .

3.1.1 Estimation via EM

Suppose all documents are independent. The log-likelihood function for $\theta = (A, w)$ is then computed as

$$\ell(\theta; X) = \sum_{i=1}^n \log P(x_i \mid \theta) \quad (119)$$

$$= \sum_{i=1}^n \log \left(\sum_{k=1}^K P(x_i, z_i = z^{(k)} \mid \theta) \right) \quad (120)$$

$$= \sum_{i=1}^n \log \left(\sum_{k=1}^K P(z_i = z^{(k)} \mid \theta) P(x_i \mid z_i = z^{(k)}, \theta) \right) \quad (121)$$

$$= \sum_{i=1}^n \log \left[\sum_{k=1}^K w(k) \prod_{j=1}^p (A_k(j))^{X(j,i)} \right], \quad (122)$$

where constants are omitted by convention. It is obvious that optimizing $\ell(\theta; X)$ does not admit a closed-form solution, which calls for effective algorithms.

¹³However, according to Ke et al. (2023), unigram models are those only model the counts of single words, neglecting word orders and word context.

¹⁴These assumptions naturally limit the model's capacity to fit a large collection of documents.

The well-known Expectation-Maximization (EM) algorithm (Dempster et al. (1977)) is a commonly used approach to fitting latent variable models. The algorithm is an iterative procedure that alternates between two steps: the Expectation (E-step) and the Maximization (M-step). In the E-step, it computes the expected value of the complete-data log-likelihood with respect to the conditional distribution of the latent variables given the observed data and current parameter estimates. In the M-step, it maximizes this expected complete-data log-likelihood with respect to the model parameters.

In the mixture of unigrams model, words and documents can be considered as *observed variables* (also referred to as *manifest variables*), while topics can be viewed as *latent variables*. The log-likelihood function for complete data (i.e., X, z_1, \dots, z_n) is computed as

$$\ell(\theta; X, z_1, \dots, z_n) = \sum_{i=1}^n \log P(x_i, z_i \mid \theta) \quad (123)$$

$$= \sum_{i=1}^n [\log P(\text{topic } z_i \mid \text{doc } i) + \log P(x_i \mid \text{topic } z_i, \text{doc } i)] \quad (124)$$

$$= \sum_{i=1}^n \left[\log P(\text{topic } z_i \mid \text{doc } i) + \sum_{j=1}^p X(j, i) \log P(\text{word } j \mid \text{topic } z_i, \text{doc } i) \right] \quad (125)$$

$$= \sum_{i=1}^n \left[\log P(\text{topic } z_i \mid \text{doc } i) + \sum_{j=1}^p X(j, i) \log P(\text{word } j \mid \text{topic } z_i) \right] \quad (126)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \left[I(z_i = z^{(k)}) \left(\log w(k) + \sum_{j=1}^p X(j, i) \log A_k(j) \right) \right]. \quad (127)$$

The **E-step** begins by computing the posterior distribution of z_i after observing x_i :

$$\begin{aligned} P(z_i = z^{(k)} \mid x_i, \theta) &= \frac{P(x_i \mid z_i = z^{(k)}) P(z_i = z^{(k)})}{\sum_{k=1}^K P(x_i \mid z_i = z^{(k)}) P(z_i = z^{(k)})} \\ &= \frac{\text{Multinomial}(x_i \mid N_i, A_k) w(k)}{\sum_{l=1}^K \text{Multinomial}(x_i \mid N_i, A_l) w(l)}, \end{aligned} \quad (128)$$

which follows from the well-known Bayes formula. Then, the so-called Q function is given by

$$Q(\theta; \theta^{(t)}) = \mathbb{E} \left[\ell(\theta; X, z_1, \dots, z_n) \mid X, \theta^{(t)} \right] \quad (129)$$

$$= \sum_{i=1}^n \sum_{k=1}^K \left[P(z_i = z^{(k)} \mid x_i, \theta^{(t)}) \left(\log w(k) + \sum_{j=1}^p X(j, i) \log A_k(j) \right) \right], \quad (130)$$

where $P(z_i = z^{(k)} \mid x_i, \theta^{(t)})$ should be computed using Eq. (128), i.e.,

$$P(z_i = z^{(k)} \mid x_i, \theta^{(t)}) = \frac{\text{Multinomial}(x_i \mid N_i, A_k^{(t)}) w(k)^{(t)}}{\sum_{l=1}^K \text{Multinomial}(x_i \mid N_i, A_l^{(t)}) w(l)^{(t)}}. \quad (131)$$

The **M-step** attempts to maximize $Q(\theta; \theta^{(t)})$ under the following constraints:

$$\sum_{j=1}^p A_k(j) = 1, \quad 1 \leq k \leq K, \quad (132)$$

$$\sum_{k=1}^K w(k) = 1, \quad (133)$$

To solve this constrained optimization problem, we introduce Lagrange multipliers $\{\tau_k\}_{k=1}^K$ and ρ , and write the Lagrangian function as

$$\mathcal{L}(\theta, \{\tau_k\}_{k=1}^K, \rho; \theta^{(t)}) = Q(\theta; \theta^{(t)}) + \sum_{k=1}^K \tau_k \left(1 - \sum_{j=1}^p A_k(j) \right) + \rho \left(1 - \sum_{k=1}^K w(k) \right). \quad (134)$$

Setting the derivatives (w.r.t. $A_k(j)$ and $w(k)$) to zero gives rise to

$$\sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)}) X(j, i) \frac{1}{A_k(j)} - \tau_k = 0, \quad 1 \leq j \leq p, 1 \leq k \leq K, \quad (135)$$

$$\sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)}) \frac{1}{w(k)} - \rho = 0, \quad 1 \leq k \leq K. \quad (136)$$

Thus, the parameter updating formula in the M-step is given by

$$A_k(j)^{(t+1)} = \frac{\sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)}) X(j, i)}{\sum_{j=1}^p \sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)}) X(j, i)}, \quad 1 \leq j \leq p, 1 \leq k \leq K, \quad (137)$$

$$w(k)^{(t+1)} = \frac{\sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)})}{\sum_{k=1}^K \sum_{i=1}^n P(z_i = z^{(k)} | x_i, \theta^{(t)})}, \quad 1 \leq k \leq K. \quad (138)$$

3.2 Hofmann's pLSI Model

[Hofmann \(1999\)](#) proposed the well-known probabilistic Latent Semantic Indexing (pLSI) model, in which each document is allowed to associated with multiple topics. We next review Hofmann's pLSI model in a slightly modified fashion.

Recall that the sets of all documents, topics, and words are denoted by $\{d^{(1)}, \dots, d^{(n)}\}$, $\{z^{(1)}, \dots, z^{(K)}\}$, and $\{w^{(1)}, \dots, w^{(p)}\}$, respectively. For each document $d^{(i)}$, the word generation process can be described as follows:

- Choose a topic z according to $P(z = z^{(k)} | d^{(i)}) = P(z^{(k)} | d^{(i)}), 1 \leq k \leq K$.
- Given the topic z , choose a word w according to $P(w = w^{(j)} | z) = P(w^{(j)} | z), 1 \leq j \leq p$.
- Repeat the above steps N_i times, where N_i is the number of words in document $d^{(i)}$.

We remark that Hofmann's pLSI model essentially assumes conditional independence between a word and a document given the topic.

The parameters to be estimated are

$$\theta = \left\{ P(z^{(k)} | d^{(i)}), P(w^{(j)} | z^{(k)}) \mid 1 \leq i \leq n, 1 \leq j \leq p, 1 \leq k \leq K \right\}. \quad (139)$$

which is of order $\mathcal{O}(K(n+p))$. Suppose all documents are independent. The log-likelihood function for θ is then computed as

$$\ell(\theta; X) = \sum_{i=1}^n \log P(x_i | \theta) \quad (140)$$

$$= \sum_{i=1}^n \sum_{j=1}^p X(j, i) \log P(w^{(j)} | d^{(i)}, \theta) \quad (141)$$

$$= \sum_{i=1}^n \sum_{j=1}^p X(j, i) \log \left[\sum_{k=1}^K P(w^{(j)}, z^{(k)} | d^{(i)}, \theta) \right] \quad (142)$$

$$= \sum_{i=1}^n \sum_{j=1}^p X(j, i) \log \left[\sum_{k=1}^K P(w^{(j)} | z^{(k)}, \theta) P(z^{(k)} | d^{(i)}, \theta) \right]. \quad (143)$$

3.2.1 An Equivalent Formulation

Recently, [Ke et al. \(2023\)](#) and [Ke and Wang \(2024\)](#) introduced an alternative formulation of pLSI, which is essentially equivalent to the original model and is detailed as follows. Specifically, they assume that

$$x_i \sim \text{Multinomial}(N_i, \Omega_i), \quad 1 \leq i \leq n, \quad (144)$$

where $\Omega_i \in \mathbb{R}^p$ is the word weight vector for document i . For each document i , they require that the word weight vector Ω_i admits the following decomposition:

$$\Omega_i = \sum_{k=1}^K w_i(k) A_k, \quad (145)$$

where $w_i = (w_i(1), \dots, w_i(p))^\top \in \mathbb{R}^p$ combines the weights of document i on different topics into a vector.

They further write the model in the matrix form. They refer to $A = [A_1, \dots, A_K] \in \mathbb{R}^{p \times K}$ and $W = [w_1, \dots, w_n] \in \mathbb{R}^{K \times n}$ as the *topic matrix* and the *topic weight matrix*, respectively. It follows immediately that

$$\Omega = AW, \quad (146)$$

where $\Omega = [\Omega_1, \dots, \Omega_n] \in \mathbb{R}^{p \times n}$.

Why Model (144)-(146) Is Equivalent to pLSI. Suppose all documents are independent. The log-likelihood function for $\theta = (A, W)$ is then computed as

$$\ell(\theta; X) = \sum_{i=1}^n \sum_{j=1}^p X(j, i) \log \Omega_i(j) \quad (147)$$

$$= \sum_{i=1}^n \sum_{j=1}^p X(j, i) \log \left[\sum_{k=1}^K A_k(j) w_i(k) \right], \quad (148)$$

which results from Eqs. (144) and (145). This log-likelihood function is equivalent to that of the pLSI model (Eq. (143)) under the reparameterization:

$$A_k(j) = P(w^{(j)} | z^{(k)}, \theta), \quad 1 \leq j \leq p, 1 \leq k \leq K, \quad (149)$$

$$w_i(k) = P(z^{(k)} | d^{(i)}, \theta), \quad 1 \leq i \leq n, 1 \leq k \leq K. \quad (150)$$

Anchor Words. When a certain word appears only when a specific topic is being discussed, we may regard this word as an *anchor word* (Arora et al. (2012)) of that topic. For example, the phrase (word) “Feixi Late Night Canteen” is used only when discussing my home institution (topic), the University of Science and Technology of China¹⁵. Rigorously, we call word j an anchor word of topic k if $A_k(j) \neq 0$ and $A_l(j) = 0$ for all $l \neq k$. This implies that when word j appears in a document, we can immediately infer that this document covers topic k (though it may also cover other topics). A practical use of anchor words is that they allow us to interpret each estimated topic and subsequently assign an appropriate label by its anchor words.

Identifiability Issue. To make Model (144)-(146) identifiable (i.e., for a given Ω , there is a unique pair (A, W) such that $\Omega = AW$ holds), we may require the *anchor-word condition* (i.e., each topic has at least one anchor word). According to Donoho and Stodden (2003) and Ke and Wang (2024), this is almost the necessary condition for identifiability of Model (144)-(146).

3.3 Latent Dirichlet Allocation

The latent Dirichlet allocation (LDA) model (Blei et al. (2003)) can be viewed as a Bayesian version of the pLSI model. In the LDA model, we assume that the topic weight vectors w_1, \dots, w_n are i.i.d. drawn from a Dirichlet distribution with parameter $\alpha = (\alpha_1, \dots, \alpha_K)$, i.e., $w_1, \dots, w_n \stackrel{\text{i.i.d.}}{\sim} \text{Dir}(\alpha)$.

Dirichlet Distribution. A random vector $\gamma = (\gamma_1, \dots, \gamma_K)^\top \in \mathbb{R}^K$ is said to be Dirichlet distributed with parameter $\alpha = (\alpha_1, \dots, \alpha_K)$, i.e., $\gamma \sim \text{Dir}(\alpha)$, if it takes value in a $K - 1$ -simplex (i.e., $\gamma_j \geq 0$ for all $1 \leq j \leq K$ and $\sum_{j=1}^K \gamma_j = 1$) and its probability density function is given by

$$p(\gamma | \alpha) = \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \gamma_i^{\alpha_i-1}, \quad (151)$$

where $\Gamma(\cdot)$ is the gamma function.

Specifically, for each document i , the generative process of LDA can be described as:

- Choose a topic distribution $w_i \sim \text{Dir}(\alpha)$.
- Repeat the following steps N_i times, where N_i is the number of words in document i :
 - Choose a topic z according to $P(z = z^{(k)} | w_i) = w_i(k), 1 \leq k \leq K$.
 - Suppose the chosen topic is $z^{(k)}$, choose a word according to $P(\text{word} = w^{(j)} | z = z^{(k)}) = A_k(j), 1 \leq j \leq p$.

Here $A_k = (A_k(1), \dots, A_k(p))^\top \in \mathbb{R}^p$ is the population word frequency vector for topic $z^{(k)}$, $1 \leq k \leq K$. We refer to $A = [A_1, \dots, A_K] \in \mathbb{R}^{p \times K}$ and $W = [w_1, \dots, w_n] \in \mathbb{R}^{K \times n}$ as the topic matrix and the topic weight matrix, respectively. Our notations here are consistent with earlier sections.

¹⁵This is a joke. Please do not take it seriously.

Given the parameters A and α , the joint distribution of x_i and w_i is given by

$$P(x_i, w_i \mid A, \alpha) = P(w_i \mid \alpha)P(x_i \mid w_i, A) \quad (152)$$

$$= \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(P(\text{word} = w^{(j)} \mid w_i, A) \right)^{X(j,i)} \quad (153)$$

$$= \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K P(\text{word} = w^{(j)}, \text{topic} = z^{(k)} \mid w_i, A) \right)^{X(j,i)} \quad (154)$$

$$= \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K P(\text{topic} = z^{(k)} \mid w_i) P(\text{word} = w^{(j)} \mid \text{topic} = z^{(k)}, A) \right)^{X(j,i)} \quad (155)$$

$$= \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K w_i(k) A_k(j) \right)^{X(j,i)}, \quad (156)$$

where constant factors (w.r.t. A and α) are omitted. Integrating over w_i gives rise to

$$P(x_i \mid A, \alpha) = \int P(x_i, w_i \mid A, \alpha) dw_i \quad (157)$$

$$= \int \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K w_i(k) A_k(j) \right)^{X(j,i)} dw_i. \quad (158)$$

Suppose all documents are independent. We take the product of the marginal probabilities of single documents and finally obtain

$$P(X \mid A, \alpha) = \prod_{i=1}^n P(x_i \mid A, \alpha) \quad (159)$$

$$= \prod_{i=1}^n \int \text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K w_i(k) A_k(j) \right)^{X(j,i)} dw_i. \quad (160)$$

The LDA model can be fitted (i.e., maximizing $P(X \mid A, \alpha)$ w.r.t (A, α)) via either the variational EM algorithm or Gibbs sampling ([Porteous et al. \(2008\)](#)). In this note, we only present the variational EM algorithm.

3.3.1 Estimation via Variational Inference

Suppose Q is some probability density function for W . The starting point of variational EM algorithm is the following derivation:

$$\log P(X | A, \alpha) = \log P(X | A, \alpha) \int Q(W) dW \quad (161)$$

$$= \int \log P(X | A, \alpha) Q(W) dW \quad (162)$$

$$= \mathbb{E}_{Q(W)} [\log P(X | A, \alpha)] \quad (163)$$

$$= \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha)}{P(W | X, A, \alpha)} \right] \quad (164)$$

$$= \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha) Q(W)}{P(W | X, A, \alpha) Q(W)} \right] \quad (165)$$

$$= \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha)}{Q(W)} \right] + \mathbb{E}_{Q(W)} \left[\log \frac{Q(W)}{P(W | X, A, \alpha)} \right] \quad (166)$$

$$= \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha)}{Q(W)} \right] + D_{\text{KL}}(Q(W) \| P(W | X, A, \alpha)) \quad (167)$$

$$\geq \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha)}{Q(W)} \right], \quad (168)$$

where the last line follows from the fact that the Kullback-Leibler (KL) divergence is always non-negative (Kullback and Leibler (1951)). The term in the last line is often referred to as the Evidence Lower BOund (ELBO), i.e.,

$$\text{ELBO}(Q, A, \alpha) = \mathbb{E}_{Q(W)} \left[\log \frac{P(X, W | A, \alpha)}{Q(W)} \right] \quad (169)$$

$$= \mathbb{E}_{Q(W)} [\log P(X, W | A, \alpha)] - \mathbb{E}_{Q(W)} [\log Q(W)]. \quad (170)$$

We can clearly observe from Eq. (167) that the log-likelihood $\log P(X | A, \alpha)$ is equal to the ELBO term $\text{ELBO}(Q, A, \alpha)$ plus the KL divergence term $D_{\text{KL}}(Q(W) \| P(W | X, A, \alpha))$. This observation motivates us to maximize $\text{ELBO}(Q, A, \alpha)$, as an alternative to directly maximizing $\log P(X | A, \alpha)$, which is the basic idea behind variational inference (Jordan et al. (1999); Wainwright et al. (2008)).

Formally, the variational EM algorithm can be stated as follows:

- **E-step:** Fix (A, α) , and solve for $Q = \arg \max_Q \text{ELBO}(Q, A, \alpha)$.
- **M-step:** Fix Q , and solve for $(A, \alpha) = \arg \max_{(A, \alpha)} \text{ELBO}(Q, A, \alpha)$.

These two steps are repeated until $\text{ELBO}(Q, A, \alpha)$ converges.

We introduce the so-called *mean-field* family, i.e.,

$$Q(W) = \prod_{i=1}^n q_i(w_i). \quad (171)$$

In the LDA model, we may simplify the computation by requiring

$$q_i(w_i) = \text{Dir}(w_i | \gamma_i), \quad 1 \leq i \leq n. \quad (172)$$

Here each $\gamma_i = (\gamma_i(1), \dots, \gamma_i(K))^\top$ is a K -dimensional vector with $\gamma_i(k) > 0$, $1 \leq k \leq K$.

Equipped with all these assumptions and tools, we write

$$\text{ELBO}(Q, A, \alpha) = \mathbb{E}_{Q(W)} [\log P(X, W \mid A, \alpha)] - \mathbb{E}_{Q(W)} [\log Q(W)] \quad (173)$$

$$= \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} [\log P(x_i, w_i \mid A, \alpha)] - \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} [\log q_i(w_i)] \quad (174)$$

$$= \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} \left\{ \log \left[\text{Dir}(w_i \mid \alpha) \prod_{j=1}^p \left(\sum_{k=1}^K w_i(k) A_k(j) \right)^{X(j,i)} \right] \right\} \quad (175)$$

$$- \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} [\log \text{Dir}(w_i \mid \gamma_i)]. \quad (176)$$

$$= \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} \left[\log \text{Dir}(w_i \mid \alpha) + \sum_{j=1}^p X(j, i) \log \left(\sum_{k=1}^K w_i(k) A_k(j) \right) \right] \quad (177)$$

$$- \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} [\log \text{Dir}(w_i \mid \gamma_i)]. \quad (178)$$

$$= \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} \left[\log \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K w_i(k)^{\alpha_k-1} + \sum_{j=1}^p X(j, i) \log \left(\sum_{k=1}^K w_i(k) A_k(j) \right) \right] \quad (179)$$

$$- \sum_{i=1}^n \mathbb{E}_{q_i(w_i)} \left[\log \frac{\Gamma(\sum_{k=1}^K \gamma_i(k))}{\prod_{k=1}^K \Gamma(\gamma_i(k))} \prod_{k=1}^K w_i(k)^{\gamma_i(k)-1} \right]. \quad (180)$$

4 Covariance Estimation

Let the matrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times p}$ denote n i.i.d. observations on a system of p random variables with covariance matrix Σ . The sample covariance matrix is given by

$$\mathbf{S}_n = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}_n) (\mathbf{x}_i - \bar{\mathbf{x}}_n)^\top \quad (181)$$

$$= \frac{1}{n} (\mathbf{X} - \mathbf{P}_{\mathbf{1}_n} \mathbf{X})^\top (\mathbf{X} - \mathbf{P}_{\mathbf{1}_n} \mathbf{X}), \quad (182)$$

where $\bar{\mathbf{x}}_n = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$ is the sample mean vector, and $\mathbf{P}_{\mathbf{1}_n} = \mathbf{1}_n (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$ is the projection matrix onto the span of $\mathbf{1}_n = (1, \dots, 1)^\top \in \mathbb{R}^n$.

In a high-dimensional setting with $p > n$, the rank of \mathbf{S}_n is strictly less than p :

$$\text{rank}(\mathbf{S}_n) = \text{rank}(\mathbf{X} - \mathbf{P}_{\mathbf{1}_n} \mathbf{X}) \quad (183)$$

$$= \text{rank}((\mathbf{I}_n - \mathbf{P}_{\mathbf{1}_n}) \mathbf{X}) \quad (184)$$

$$\leq \text{rank}(\mathbf{I}_n - \mathbf{P}_{\mathbf{1}_n}) \quad (185)$$

$$= n - 1 < p, \quad (186)$$

which implies that the sample covariance matrix \mathbf{S}_n is not invertible. Moreover, we can examine

the Frobenius distance between \mathbf{S}_n and $\mathbf{\Sigma}$. Specifically, we have¹⁶

$$\mathbb{E} \left[\|\mathbf{S}_n - \mathbf{\Sigma}\|_{\text{F}}^2 \right] = \mathbb{E} \left[\sum_{i=1}^p \sum_{j=1}^p (S_n(i, j) - \Sigma(i, j))^2 \right] \quad (187)$$

$$= \sum_{i=1}^p \sum_{j=1}^p \mathbb{E} [S_n(i, j) - \Sigma(i, j)]^2 \quad (188)$$

$$= \sum_{i=1}^p \sum_{j=1}^p \text{Var} [S_n(i, j)] \quad (189)$$

$$= \sum_{i=1}^p \sum_{j=1}^p \mathcal{O}\left(\frac{1}{n}\right) \quad (190)$$

$$= \mathcal{O}\left(\frac{p^2}{n}\right). \quad (191)$$

Thus, the sample covariance matrix does not converge to $\mathbf{\Sigma}$ in the Frobenius norm when $p > n$.

In this section, we review several classical estimators that are both invertible and consistent in high-dimensional settings.

4.1 A Linear Shrinkage Estimator

Ledoit and Wolf (2004) proposed an estimator which is essentially a linear combination of the sample covariance matrix \mathbf{S}_n with the identity matrix \mathbf{I}_p . The intuition arises from the following optimization problem:

$$\min_{\rho_1, \rho_2} \mathbb{E} \left[\left\| \hat{\mathbf{\Sigma}} - \mathbf{\Sigma} \right\|_{\text{F}}^2 \right], \quad \text{s.t. } \hat{\mathbf{\Sigma}} = \rho_1 \mathbf{I}_p + \rho_2 \mathbf{S}_n, \quad (192)$$

where ρ_1 and ρ_2 are non-random coefficients. Interestingly, this optimization problem admits a closed-form solution, which we will present in the sequel.

We first introduce the inner product $\langle \cdot, \cdot \rangle_{\text{F}}$ induced by $\|\cdot\|_{\text{F}}$, i.e.,

$$\langle \mathbf{A}, \mathbf{B} \rangle_{\text{F}} = \frac{1}{4} \left(\|\mathbf{A} + \mathbf{B}\|_{\text{F}}^2 - \|\mathbf{A} - \mathbf{B}\|_{\text{F}}^2 \right) \quad (193)$$

$$= \frac{1}{4} \left[\text{tr} \left((\mathbf{A} + \mathbf{B})(\mathbf{A} + \mathbf{B})^{\top} \right) - \text{tr} \left((\mathbf{A} - \mathbf{B})(\mathbf{A} - \mathbf{B})^{\top} \right) \right] \quad (194)$$

$$= \text{tr} \left(\mathbf{A} \mathbf{B}^{\top} \right). \quad (195)$$

Four scalars play a central role in the analysis:

$$\mu = \langle \mathbf{\Sigma}, \mathbf{I}_p \rangle_{\text{F}} = \text{tr}(\mathbf{\Sigma}), \quad (196)$$

$$\alpha^2 = \|\mathbf{\Sigma} - \mu \mathbf{I}_p\|_{\text{F}}^2, \quad (197)$$

$$\beta^2 = \mathbb{E} \left[\|\mathbf{S}_n - \mathbf{\Sigma}\|_{\text{F}}^2 \right], \quad (198)$$

$$\delta^2 = \mathbb{E} \left[\|\mathbf{S}_n - \mu \mathbf{I}_p\|_{\text{F}}^2 \right]. \quad (199)$$

¹⁶The same result is provided by Fan et al. (2008).

It follows immediately that

$$\delta^2 = \mathbb{E} \left[\|\mathbf{S}_n - \mu \mathbf{I}_p\|_F^2 \right] \quad (200)$$

$$= \mathbb{E} \left[\|(\mathbf{S}_n - \boldsymbol{\Sigma}) + (\boldsymbol{\Sigma} - \mu \mathbf{I}_p)\|_F^2 \right] \quad (201)$$

$$= \mathbb{E} \left[\|\mathbf{S}_n - \boldsymbol{\Sigma}\|_F^2 \right] + 2\mathbb{E} [\langle \mathbf{S}_n - \boldsymbol{\Sigma}, \boldsymbol{\Sigma} - \mu \mathbf{I}_p \rangle_F] + \|\boldsymbol{\Sigma} - \mu \mathbf{I}_p\|_F^2 \quad (202)$$

$$= \beta^2 + 0 + \alpha^2 = \alpha^2 + \beta^2. \quad (203)$$

We are now ready to solve the optimization problem (192). Specifically, we have

$$\min_{\rho_1, \rho_2} \mathbb{E} \left[\left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma} \right\|_F^2 \right] = \min_{\rho_1, \rho_2} \mathbb{E} \left[\|(\rho_1 \mathbf{I}_p + \rho_2 \mathbf{S}_n) - \boldsymbol{\Sigma}\|_F^2 \right] \quad (204)$$

$$= \min_{\rho, v} \mathbb{E} \left[\|(\rho v \mathbf{I}_p + (1 - \rho) \mathbf{S}_n) - \boldsymbol{\Sigma}\|_F^2 \right] \quad (\text{A change of variables}) \quad (205)$$

$$= \min_{\rho, v} \left\{ \rho^2 \|\boldsymbol{\Sigma} - v \mathbf{I}_p\|_F^2 + (1 - \rho)^2 \mathbb{E} \left[\|\mathbf{S}_n - \boldsymbol{\Sigma}\|_F^2 \right] \right\} \quad (206)$$

$$= \min_{\rho, v} \left\{ \rho^2 \|\boldsymbol{\Sigma} - v \mathbf{I}_p\|_F^2 + (1 - \rho)^2 \beta^2 \right\} \quad (207)$$

$$= \min_{\rho} \left\{ \rho^2 \|\boldsymbol{\Sigma} - \mu \mathbf{I}_p\|_F^2 + (1 - \rho)^2 \beta^2 \right\} \quad (v = \langle \boldsymbol{\Sigma}, \mathbf{I}_p \rangle_F = \mu) \quad (208)$$

$$= \min_{\rho} \left\{ \rho^2 \alpha^2 + (1 - \rho)^2 \beta^2 \right\} \quad (209)$$

$$= \left(\frac{\beta^2}{\delta^2} \right)^2 \alpha^2 + \left(1 - \frac{\beta^2}{\delta^2} \right)^2 \beta^2 \quad \left(\rho = \frac{\beta^2}{\alpha^2 + \beta^2} = \frac{\beta^2}{\delta^2} \right) \quad (210)$$

$$= \frac{\alpha^2 \beta^2}{\delta^2}. \quad (\delta^2 = \alpha^2 + \beta^2) \quad (211)$$

Thus, the optimal $\hat{\boldsymbol{\Sigma}}$ is given by

$$\hat{\boldsymbol{\Sigma}} = \frac{\beta^2}{\delta^2} \mu \mathbf{I}_p + \frac{\alpha^2}{\delta^2} \mathbf{S}_n, \quad (212)$$

where $\frac{\beta^2}{\delta^2} \in [0, 1]$ can be interpreted as the shrinkage intensity placed on the shrinkage target $\mu \mathbf{I}_p$. The percentage relative improvement in average loss over the sample covariance matrix is equal to

$$\frac{\mathbb{E} \left[\|\mathbf{S}_n - \boldsymbol{\Sigma}\|_F^2 \right] - \mathbb{E} \left[\left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma} \right\|_F^2 \right]}{\mathbb{E} \left[\|\mathbf{S}_n - \boldsymbol{\Sigma}\|_F^2 \right]} = \frac{\beta^2 - \frac{\alpha^2 \beta^2}{\delta^2}}{\beta^2} = \frac{\beta^2}{\delta^2}. \quad (213)$$

Intuitively, if \mathbf{S}_n is relatively accurate ($\beta^2/\delta^2 \approx 0$), then we should not shrink it too much; if \mathbf{S}_n is relatively accurate ($\beta^2/\delta^2 \approx 1$), we expect to gain a lot by applying the shrinkage.

In practice, we never have access to the ground-truth μ , α , β , and δ . [Ledoit and Wolf \(2004\)](#) suggested we estimate these quantities by

$$\hat{\mu} = \langle \mathbf{S}_n, \mathbf{I}_p \rangle_F = \text{tr}(\mathbf{S}_n), \quad (214)$$

$$\hat{\delta}^2 = \|\mathbf{S}_n - \hat{\mu} \mathbf{I}_p\|_F^2, \quad (215)$$

$$\widehat{\beta^2} = \frac{1}{n} \sum_{i=1}^n \left\| (\mathbf{x}_i - \bar{\mathbf{x}}_n) (\mathbf{x}_i - \bar{\mathbf{x}}_n)^\top - \mathbf{S}_n \right\|_F^2, \quad (216)$$

$$\widehat{\alpha^2} = \hat{\delta}^2 - \widehat{\beta^2}. \quad (217)$$

Finally, the estimator is given by

$$\hat{\Sigma} = \frac{\widehat{\beta^2}}{\widehat{\delta^2}} \hat{\mu} \mathbf{I}_p + \frac{\widehat{\alpha^2}}{\widehat{\delta^2}} \mathbf{S}_n, \quad (218)$$

which is consistent (in the sense that $\mathbb{E} \left[\left\| \hat{\Sigma} - \Sigma \right\|_{\text{F}}^2 \right] \rightarrow 0$ as $n \rightarrow \infty$) under some suitable conditions¹⁷.

4.2 Generalized Thresholding

Rothman et al. (2009) proposed a class of generalized thresholding operators that combine thresholding with shrinkage, and studied its application to sample covariance estimation in high dimensions. Here we briefly review their method.

For any $\lambda \geq 0$, we define a generalized thresholding operator $s_\lambda : \mathbb{R} \rightarrow \mathbb{R}$ that satisfies the following conditions¹⁸ for all $z \in \mathbb{R}$:

- (i) **Shrinkage:** $|s_\lambda(z)| \leq |z|$;
- (ii) **Thresholding:** $s_\lambda(z) = 0$ for $|z| \leq \lambda$;
- (iii) **The Amount of Shrinkage:** $|s_\lambda(z) - z| \leq \lambda$.

It turns out that the above conditions are satisfied by a number of commonly used shrinkage/thresholding rules, including the hard-thresholding rule $s_\lambda^{\text{H}}(z) = z \cdot I(|z| \geq \lambda)$, the LASSO $s_\lambda^{\text{L}}(z) = \text{sign}(z)(|z| - \lambda)_+$ (Tibshirani (1996)), and the adaptive LASSO $s_\lambda^{\text{AL}}(z) = \text{sign}(z)(|z| - \lambda^{\eta+1}|z|^{-\eta})_+$ with $\eta \geq 1$ (Zou (2006)).

Then, we apply the operator entry-wise to the sample covariance matrix, yielding an estimator denoted by $s_\lambda(\mathbf{S}_n)$, i.e.,

$$(s_\lambda(\mathbf{S}_n))(i, j) = s_\lambda(S_n(i, j)), \quad 1 \leq i, j \leq p. \quad (219)$$

Rothman et al. (2009) proved that $s_\lambda(\mathbf{S}_n)$ is an L_2 -consistent estimator of Σ under suitable conditions¹⁹.

4.2.1 Proof of Consistency

In this section, we establish the consistency of the generalized thresholding estimator $s_\lambda(\mathbf{S}_n)$. The proof is largely benefited from the original paper Rothman et al. (2009).

We impose a condition on the marginal distribution of x_{1j} , $1 \leq j \leq p$, where x_{1j} is the j -th component of \mathbf{x}_1 (see, e.g., Bickel and Levina (2008a) and Bickel and Levina (2008b)):

$$\mathbb{E} [\exp(\lambda x_{1j}^2)] = \int_0^\infty \exp(\lambda t) dG_j(t) < \infty, \quad 0 \leq |\lambda| \leq \lambda_0, \quad (220)$$

¹⁷In the original paper Ledoit and Wolf (2004), one of the assumptions is that there exists a constant K independent of n such that $p/n \leq K$, which is pretty weak.

¹⁸In principle, it is acceptable to use different parameters λ_1 and λ_2 in (ii) and (iii). Here we set them to be the same for simplicity.

¹⁹One of the assumptions is that $\log p/n \rightarrow 0$ as $n \rightarrow \infty$, which is obviously weaker than $p^2/n \rightarrow 0$.

for some $\lambda_0 > 0$, where G_j is the cumulative distribution function of x_{1j}^2 .

The consistency result we aim to establish holds uniformly on a class of “approximately sparse” covariance matrices, which was introduced by [Bickel and Levina \(2008a\)](#):

$$\mathcal{U}_\tau(q, c_0(p), M) = \left\{ \mathbf{\Sigma} : \max_i \Sigma(i, i) \leq M, \max_i \sum_{j=1}^p |\Sigma(i, j)|^q \leq c_0(p) \right\}, \quad (221)$$

for some $0 \leq q < 1$. Note that when $q = 0$, this is a class of truly sparse matrices with

$$\mathcal{U}_\tau(0, c_0(p), M) = \left\{ \mathbf{\Sigma} : \max_i \Sigma(i, i) \leq M, \max_i \sum_{j=1}^p I(\Sigma(i, j) \neq 0) \leq c_0(p) \right\}. \quad (222)$$

According to [Bickel and Levina \(2008a\)](#), if Assumption (220) holds and $\mathbf{\Sigma} \in \mathcal{U}_\tau(q, c_0(p), M)$, the following useful results apply:

$$\max_i \sum_{j=1}^p |S_n(i, j)| I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| < \lambda) = \mathcal{O}_P \left(c_0(p) \lambda^{-q} \sqrt{\frac{\log p}{n}} + c_0(p) \lambda^{1-q} \right), \quad (223)$$

$$\max_i \sum_{j=1}^p |\Sigma(i, j)| I(|S_n(i, j)| < \lambda, |\Sigma(i, j)| \geq \lambda) = \mathcal{O}_P \left(c_0(p) \lambda^{-q} \sqrt{\frac{\log p}{n}} + c_0(p) \lambda^{1-q} \right), \quad (224)$$

$$\max_i \sum_{j=1}^p |S_n(i, j) - \Sigma(i, j)| I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) = \mathcal{O}_P \left(c_0(p) \lambda^{-q} \sqrt{\frac{\log p}{n}} \right), \quad (225)$$

which will play an important role in the proof presented below.

The proof begins with the following triangle inequality

$$\|s_\lambda(\mathbf{S}_n) - \mathbf{\Sigma}\|_2 \leq \|s_\lambda(\mathbf{\Sigma}) - \mathbf{\Sigma}\|_2 + \|s_\lambda(\mathbf{S}_n) - s_\lambda(\mathbf{\Sigma})\|_2. \quad (226)$$

For any symmetric matrix \mathbf{A} , we have (see, e.g., Proposition 9.1 in the textbook [Fan et al. \(2020\)](#))

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_1 = \|\mathbf{A}\|_\infty = \max_i \sum_j |A(i, j)|. \quad (227)$$

In light of this result, we are going to bound all the operator norms by (227).

For the first term in (226), we have

$$\sum_{j=1}^p |s_\lambda(\Sigma(i, j)) - \Sigma(i, j)| \leq \sum_{j=1}^p |\Sigma(i, j)| I(|\Sigma(i, j)| \leq \lambda) + \sum_{j=1}^p \lambda I(|\Sigma(i, j)| > \lambda) \quad (228)$$

$$= \sum_{j=1}^p |\Sigma(i, j)|^q |\Sigma(i, j)|^{1-q} I(|\Sigma(i, j)| \leq \lambda) + \sum_{j=1}^p \lambda^q \lambda^{1-q} I(|\Sigma(i, j)| > \lambda) \quad (229)$$

$$\leq \lambda^{1-q} \sum_{j=1}^p |\Sigma(i, j)|^q \quad (230)$$

$$\leq \lambda^{1-q} c_0(p), \quad (231)$$

where the first line follows from properties (ii) and (iii) of $s_\lambda(\cdot)$, and the last line follows from $\Sigma \in \mathcal{U}_\tau(q, c_0(p), M)$. Thus, the first term in (226) is bounded by $\lambda^{1-q}c_0(p)$, i.e.,

$$\|s_\lambda(\Sigma) - \Sigma\|_2 \leq \max_i \sum_{j=1}^p |s_\lambda(\Sigma(i, j)) - \Sigma(i, j)| \quad (232)$$

$$\leq \lambda^{1-q}c_0(p). \quad (233)$$

For the second term in (226), utilizing properties (i) and (ii) of $s_\lambda(\cdot)$, we have

$$|s_\lambda(S_n(i, j)) - s_\lambda(\Sigma(i, j))| \quad (234)$$

$$\leq |S_n(i, j)| I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| < \lambda) \quad (235)$$

$$+ |\Sigma(i, j)| I(|S_n(i, j)| < \lambda, |\Sigma(i, j)| \geq \lambda) \quad (236)$$

$$+ (|S_n(i, j) - \Sigma(i, j)| + |s_\lambda(S_n(i, j)) - S_n(i, j)| + |s_\lambda(\Sigma(i, j)) - \Sigma(i, j)|) \quad (237)$$

$$\cdot I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda). \quad (238)$$

The first three terms can be bounded via (223)-(225), respectively. For the fourth term, applying (iii), we have

$$\max_i \sum_{j=1}^p |s_\lambda(S_n(i, j)) - S_n(i, j)| I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) \quad (239)$$

$$\leq \max_i \sum_{j=1}^p \lambda I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) \quad (240)$$

$$= \lambda^{1-q} \max_i \sum_{j=1}^p \lambda^q I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) \quad (241)$$

$$\leq \lambda^{1-q} \max_i \sum_{j=1}^p |\Sigma(i, j)|^q I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) \quad (242)$$

$$\leq \lambda^{1-q} \max_i \sum_{j=1}^p |\Sigma(i, j)|^q \quad (243)$$

$$\leq \lambda^{1-q}c_0(p). \quad (244)$$

Similarly, for the last term, we have

$$\max_i \sum_{j=1}^p |s_\lambda(\Sigma(i, j)) - \Sigma(i, j)| I(|S_n(i, j)| \geq \lambda, |\Sigma(i, j)| \geq \lambda) \leq \lambda^{1-q}c_0(p). \quad (245)$$

Finally, by collecting the terms, we obtain

$$\|s_\lambda(\mathbf{S}_n) - \Sigma\|_2 \leq \|s_\lambda(\Sigma) - \Sigma\|_2 + \|s_\lambda(\mathbf{S}_n) - s_\lambda(\Sigma)\|_2 \quad (246)$$

$$\leq 3\lambda^{1-q}c_0(p) + 2\mathcal{O}_P\left(c_0(p)\lambda^{-q}\sqrt{\frac{\log p}{n}} + c_0(p)\lambda^{1-q}\right) + \mathcal{O}_P\left(c_0(p)\lambda^{-q}\sqrt{\frac{\log p}{n}}\right) \quad (247)$$

$$= \mathcal{O}_P\left(c_0(p)\lambda^{-q}\sqrt{\frac{\log p}{n}} + c_0(p)\lambda^{1-q}\right). \quad (248)$$

The optimal choice $\lambda = \lambda^* \asymp \sqrt{\frac{\log p}{n}}$ gives rise to the following convergence rate:

$$\|s_{\lambda^*}(\mathbf{S}_n) - \mathbf{\Sigma}\|_2 \leq \mathcal{O}_P \left(c_0(p) \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}} \right). \quad (249)$$

4.3 The POET Estimator

Suppose the samples $\{\mathbf{x}_i\}_{i=1}^n$ are independently drawn from the following orthogonal factor model:

$$\mathbf{x} = \mathbf{L}\mathbf{f} + \boldsymbol{\epsilon}, \quad (250)$$

where

- $\mathbf{L} \in \mathbb{R}^{p \times K}$ is the non-random loading matrix with $\text{rank}(\mathbf{L}) = K < p$,
- \mathbf{f} is a K -dimensional random vector of common factors with $\text{Cov}[\mathbf{f}] = \mathbf{I}_K$,
- $\boldsymbol{\epsilon}$ is a p -dimensional random noise vector, independent of \mathbf{f} , with $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}_p$.

It follows immediately that

$$\mathbf{\Sigma} = \text{Cov}[\mathbf{x}] \quad (251)$$

$$= \mathbf{L}\mathbf{L}^\top + \text{Cov}[\boldsymbol{\epsilon}], \quad (252)$$

implying that $\mathbf{\Sigma}$ should be “close” to a rank- K matrix. In fact, according to the [Eckart-Young-Mirsky low-rank approximation theorem](#), the best rank- K approximation to $\mathbf{\Sigma}$ in either the Frobenius norm or the spectral norm is given by its top- K singular value decomposition (SVD). In this context, we next review the Principal Orthogonal complement Thresholding (POET) estimator proposed by [Fan et al. \(2013\)](#).

Let $\hat{\lambda}_1, \dots, \hat{\lambda}_p$ be ordered eigenvalues of \mathbf{S}_n and $\{\hat{\boldsymbol{\xi}}_i\}_{i=1}^p$ be their corresponding eigenvectors. Let K denote the number of diverging eigenvalues. Then the sample covariance admits the following spectral decomposition:

$$\mathbf{S}_n = \sum_{i=1}^K \hat{\lambda}_i \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^\top + \hat{\mathbf{R}}_K, \quad (253)$$

where $\hat{\mathbf{R}}_K = \sum_{i=K+1}^p \hat{\lambda}_i \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^\top$ is the principal orthogonal complement. [Fan et al. \(2013\)](#) retained $\sum_{i=1}^K \hat{\lambda}_i \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^\top$ unchanged and applied entry-wise thresholding to $\hat{\mathbf{R}}_K$ via²⁰

$$\hat{R}_K^\mathcal{T}(i, j) = s_{\lambda_{ij}} \left(\hat{R}_K(i, j) \right), \quad 1 \leq i, j \leq p, \quad (254)$$

where λ_{ij} ’s are entry-dependent threshold parameters. The POET estimator is then given by

$$\hat{\mathbf{\Sigma}}_K^{\text{POET}} = \sum_{i=1}^K \hat{\lambda}_i \hat{\boldsymbol{\xi}}_i \hat{\boldsymbol{\xi}}_i^\top + \hat{\mathbf{R}}_K^\mathcal{T}. \quad (255)$$

[Fan et al. \(2013\)](#) proved that both $\hat{\mathbf{\Sigma}}_K^{\text{POET}}$ and $\left(\hat{\mathbf{\Sigma}}_K^{\text{POET}} \right)^{-1}$ are consistent under suitable conditions.

²⁰In the original paper [Fan et al. \(2013\)](#), the diagonal entries are not subject to thresholding, i.e., $\hat{R}_K^\mathcal{T}(i, i) = \hat{R}_K(i, i), 1 \leq i \leq p$. This is equivalent to setting $\lambda_{ii} = 0$ for all $1 \leq i \leq p$.

5 Deep Learning

Exposure to Deep Learning. During my junior year of college, I enrolled in a course titled Introduction to Deep Learning²¹, earning a 98 in the class. Through hands-on projects, I gained valuable experience implementing

- Convolutional neural networks,
- Recurrent neural networks,
- Graphical neural networks,
- Large language model fine-tuning.

Furthermore, during my research internship at UM, I participated in a seminar on diffusion models led by Professor [Jeffrey Regier](#) and Professor [Yang Chen](#), where we discussed recent developments in diffusion models from both theoretical and empirical perspectives.

The Success of Deep Learning. It is generally accepted that

“the price to pay for achieving low bias is high variance” ([Geman et al. \(1992\)](#)),

a principle commonly referred to as the *bias-variance tradeoff*. From a statistical perspective, deep learning models are over-parameterized, with the number of parameters (weights and biases) extremely larger than the number of samples. Such complex models typically suffer from large variances, leading to poor performance on test sets. In my view, the success of deep learning can be attributed to the following two key factors:

- **The Arrival of Big Data.** The huge sample size, which is commonly encountered in the era of big data, reduces the variance to an acceptable level. This is analogous to the fact that the variance of a kernel estimator at a given point is $\mathcal{O}(\frac{1}{nh})$, where n is the sample size and $h > 0$ is the bandwidth (see, e.g., Proposition 1.1 in [Tsybakov \(2009\)](#)).
- **Modern Computing Power.** It is typically challenging to optimize the loss function for an over-parameterized model. However, advancements in modern computing technologies (e.g., GPUs, TPUs) have made such optimization feasible.

5.1 Diffusion Models

The main purpose of the denoising diffusion probability model (DDPM; [Ho et al. \(2020\)](#)) is to generate data from a desired distribution, especially for image data. As early as 2015, [Sohl-Dickstein et al. \(2015\)](#) introduced the idea of using a Markov chain of transitions between latent states. They referred to the encoder as the forward process, and the decoder as the backward process (Figure 1).

²¹An elective course for statistics undergraduates taught by the School of Artificial Intelligence and Data Science at USTC.

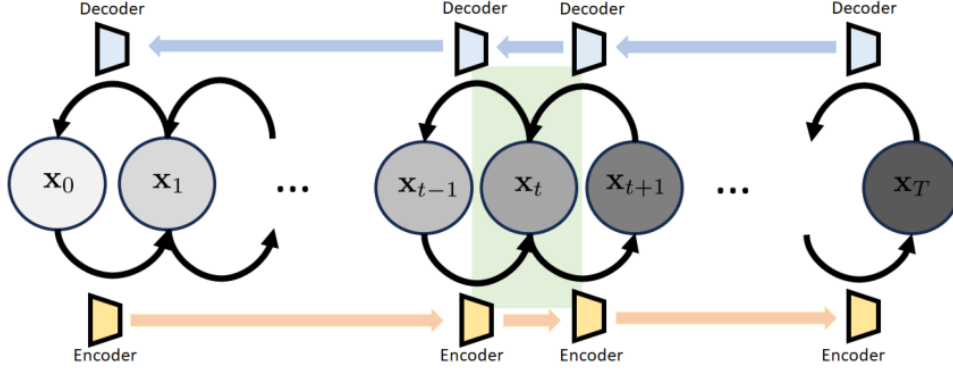


Figure 1: A variational diffusion model by Kingma et al. (2021). The input image is \mathbf{x}_0 and the white noise is \mathbf{x}_T . The intermediate states $\mathbf{x}_1, \dots, \mathbf{x}_{T-1}$ are latent variables.

Comparison with VAE. The variational autoencoder (VAE; Kingma (2013)) is an earlier tool for image generation. In a VAE, the latent variable denoted by \mathbf{z} , is typically assumed to be standard gaussian distributed, i.e., $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The input vector (typically an image) is denoted by \mathbf{x} . The central spirit of an VAE is to learn the conditional distributions $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ associated with the encoder and decoder, respectively, with these two conditional distributions parameterized by separate neural networks. The VAE is, intuitively speaking, a one-step transition process. However, DDPM **incrementally** updates the latent variables where the assembly of the whole forms the encoder-decoder structure. Other distinctions between DDPMs and VAEs include:

- **Latent Dimension.** In a DDPM, the dimensionality of all latent states is the same as that of the input \mathbf{x}_0 . In a VAE, the latent variable \mathbf{z} typically has a smaller dimensionality than the input \mathbf{x} , because we want the latent low-dimensional variable \mathbf{z} to capture the essential information required to describe \mathbf{x} .
- **Forward and Backward Processes.** In a DDPM, the forward process admits a closed-form scheme in that the conditional distribution of \mathbf{x}_t given \mathbf{x}_{t-1} , denoted by $q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1})$, is given by $\mathcal{N}(\mathbf{x}_t | \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I})$. However, in a VAE, the conditional distributions $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ are simple gaussian distributions parameterized by neural networks (which must be learned).

5.1.1 DDPM

In a DDPM, the forward process adds Gaussian noise step by step:

$$q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathbf{I}), \quad 1 \leq t \leq T, \quad (256)$$

where $0 < \alpha_t < 1$ are scalars controlling the signal-to-noise ratio at each step.

Under the Markov assumption, we can express the direct transition from \mathbf{x}_0 to \mathbf{x}_t as:

$$q_\phi(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \quad 1 \leq t \leq T, \quad (257)$$

where $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. The role of α_t 's is to ensure that $q_\phi(\mathbf{x}_T|\mathbf{x}_0) \rightarrow \mathcal{N}(\mathbf{0}, \mathbf{I})$ as $T \rightarrow \infty$. For example, a naive choice of $\alpha_1 = \dots = \alpha_T$ satisfies this requirement.

For the backward process, we choose $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ to be a Gaussian:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}\left(\mathbf{x}_{t-1} \mid \underbrace{\boldsymbol{\mu}_\theta(\mathbf{x}_t)}_{\text{neural network}}, \sigma_q^2(t)\mathbf{I}\right), \quad 1 \leq t \leq T. \quad (258)$$

Note that the neural network $\boldsymbol{\mu}_\theta(\cdot)$ is shared across all steps²². This choice is enlightened by the following fact:

$$q_\phi(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0), \sigma_q^2(t)\mathbf{I}), \quad 1 \leq t \leq T, \quad (259)$$

where

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0, \quad 1 \leq t \leq T, \quad (260)$$

$$\sigma_q^2(t) = \frac{(1 - \alpha_t)(1 - \sqrt{\bar{\alpha}_{t-1}})}{1 - \bar{\alpha}_t}, \quad 1 \leq t \leq T. \quad (261)$$

The ELBO for the above DDPM, ignoring constants, is given by (see, e.g., Theorem 2.6 in Chan (2024))

$$\text{ELBO}(\boldsymbol{\theta}) = - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \left\| \boldsymbol{\mu}_\theta(\mathbf{x}_t) - \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) \right\|_2^2 \right]. \quad (262)$$

Since $\boldsymbol{\mu}_\theta(\cdot)$ is our *design*, we rewrite it in a more convenient form:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \underbrace{\mathbf{x}_\theta(\mathbf{x}_t)}_{\text{another neural network}}, \quad 1 \leq t \leq T, \quad (263)$$

which finally gives rise to

$$\text{ELBO}(\boldsymbol{\theta}) = - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1 - \alpha_t)^2 \bar{\alpha}_{t-1}}{(1 - \bar{\alpha}_t)^2} \left\| \mathbf{x}_\theta(\mathbf{x}_t) - \mathbf{x}_0 \right\|_2^2 \right]. \quad (264)$$

The training algorithm for a DDPM is then given by:

- For every image \mathbf{x}_0 in the dataset, repeat the following steps until convergence.
- Choose $t \sim \text{Uniform}\{1, \dots, T\}$.
- Draw $\mathbf{x}_t^{(m)}$'s according to $\mathbf{x}_t^{(m)} \sim \mathcal{N}(\mathbf{x}_t \mid \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$, $1 \leq m \leq M$.
- Take gradient descent step on $\nabla_{\boldsymbol{\theta}} \left\{ \frac{1}{M} \sum_{m=1}^M \left\| \mathbf{x}_\theta(\mathbf{x}_t^{(m)}) - \mathbf{x}_0 \right\|_2^2 \right\}$.

With the trained denoiser $\mathbf{x}_\theta(\cdot)$, the inference²³ algorithm is given by:

- Draw $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- For $t = T, T-1, \dots, 1$, update according to:

$$\mathbf{x}_{t-1} = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_\theta(\mathbf{x}_t) + \sigma_q(t)\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (265)$$

²²Here we mean the time step t should also be an input to the network $\boldsymbol{\mu}_\theta(\cdot)$.

²³In the context of generative models, the term “inference” refers to “sampling”, which is different from its conventional meaning in the statistical literature.

5.1.2 A Score-Based Interpretation

Song et al. (2020) introduced an interesting interpretation of DDPMs, which we briefly discuss here.

Given the conditional distribution $q_\phi(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$ and a sample \mathbf{x}_t , Tweedie’s Formula (Efron (2011)) calculates the posterior mean of \mathbf{x}_t by

$$\mathbb{E}[\sqrt{\bar{\alpha}_t}\mathbf{x}_0 | \mathbf{x}_t] = \mathbf{x}_t + (1 - \bar{\alpha}_t)\mathbf{I}\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \quad (266)$$

$$= \mathbf{x}_t + (1 - \bar{\alpha}_t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t), \quad (267)$$

where $p_t(\cdot)$ is the marginal probability density function of \mathbf{x}_t , i.e.,

$$p_t(\mathbf{x}_t) = \int q_\phi(\mathbf{x}_t|\mathbf{x}_0) p_0(\mathbf{x}_0) d\mathbf{x}_0. \quad (268)$$

Here $p_0(\cdot)$ denotes the probability density function of the input data \mathbf{x}_0 . The term $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)$ is commonly referred to as the Stein’s score function of \mathbf{x}_t . Below we will demonstrate that a DDPM essentially learns this score function.

The above derivation suggests that

$$\mathbf{x}_0 \stackrel{\text{in some sense}}{\approx} \frac{1}{\sqrt{\bar{\alpha}_t}} (\mathbf{x}_t + (1 - \bar{\alpha}_t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)). \quad (269)$$

Plugging this approximation into our ground-truth transition mean $\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)$ yields

$$\boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0) = \frac{(1 - \bar{\alpha}_{t-1})\sqrt{\alpha_t}}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{(1 - \alpha_t)\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_t} \mathbf{x}_0, \quad (270)$$

$$\approx \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \quad (271)$$

This motivates us to redefine $\boldsymbol{\mu}_\theta(\cdot)$ as

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t) = \frac{1}{\sqrt{\alpha_t}} \mathbf{x}_t + \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \underbrace{\mathbf{s}_\theta(\mathbf{x}_t, t)}_{\text{neural network}}, \quad (272)$$

and rewrite the ELBO as

$$\text{ELBO}(\boldsymbol{\theta}) = - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \|\boldsymbol{\mu}_\theta(\mathbf{x}_t) - \boldsymbol{\mu}_q(\mathbf{x}_t, \mathbf{x}_0)\|_2^2 \right] \quad (273)$$

$$\approx - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{\alpha_t} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \right]. \quad (274)$$

Since we are more interested in the average effect over the entire dataset rather than a single

sample, we consider

$$- \sum_{t=1}^T \mathbb{E}_{p_0(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \right] \quad (275)$$

$$= - \sum_{t=1}^T \mathbb{E}_{p_t(\mathbf{x}_t)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)\|_2^2 \right] \quad (\mathbf{x}_0 \text{ is not inside the expectation}) \quad (276)$$

$$= - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_\phi(\mathbf{x}_t|\mathbf{x}_0)\|_2^2 \right] \quad (\text{See, e.g., Section 4 in Vincent (2011)}) \quad (277)$$

$$= - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t, \mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) - \left(-\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \right) \right\|_2^2 \right] \quad (278)$$

$$= - \sum_{t=1}^T \mathbb{E}_{p_0(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\frac{1}{\sigma_q^2(t)} \frac{(1-\alpha_t)^2}{\alpha_t} \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) - \left(-\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \right) \right\|_2^2 \right]. \quad (279)$$

The training algorithm is then given by:

- For every image \mathbf{x}_0 in the dataset, repeat the following steps until convergence.
- Choose $t \sim \text{Uniform}\{1, \dots, T\}$.
- Draw $\mathbf{x}_t \sim \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$.
- Take gradient descent step on $\nabla_{\theta} \left\| \mathbf{s}_\theta(\mathbf{x}_t, t) - \left(-\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0}{1 - \bar{\alpha}_t} \right) \right\|_2^2$.

With the trained score estimator $\mathbf{s}_\theta(\cdot, \cdot)$, we generate samples via Langevin MCMC (see, e.g., Song and Ermon (2019) and Section 3 in Chan (2024)):

- Draw $\mathbf{x}_T^{(0)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- For $t = T, T-1, \dots, 1$, update according to (here $\tau > 0$ is the step size)

$$\mathbf{x}_t^{(m)} = \mathbf{x}_t^{(m-1)} + \tau \mathbf{s}_\theta(\mathbf{x}_t^{(m-1)}, t) + \sqrt{2\tau} \boldsymbol{\epsilon}_t^{(m)}, \quad \boldsymbol{\epsilon}_t^{(m)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad 1 \leq m \leq M, \quad (280)$$

$$\mathbf{x}_{t-1}^{(0)} = \mathbf{x}_t^{(M)}. \quad (281)$$

- Finally, the generated sample $\mathbf{x}_0^{(0)}$ can be viewed as a draw from $p_0(\cdot)$.

5.1.3 Theoretical Foundation

Despite the huge empirical success, the theory of diffusion models is still in its infancy. Vincent (2011) revealed that the training of denoising networks essentially learns the Stein's score function denoted by $\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$, where $p(\cdot)$ is the probability density function. Therefore, diffusion models fall into the category of score-based generative models (Song et al. (2020)). A recent and highly cited (> 100 on Google Scholar) paper Chen et al. (2023), utilized this insight and demonstrated that

- The Stein’s score function can be approximated by a neural network, when the input \mathbf{x} admits a low-dimensional linear representation $\mathbf{x} = A\mathbf{z}$.
- Distribution estimation guarantee (Theorem 3 therein) can be derived using the learned score estimator.

Briefly speaking, [Chen et al. \(2023\)](#) provides insights into distribution recovery from a score-based sampling perspective. Specifically, the sample complexities depend on the intrinsic dimension of \mathbf{x} (i.e., the dimensionality of \mathbf{z}), and are free from the curse of ambient dimensionality.

5.2 Large Language Models

The entire training process of a Large Language Model (LLM) typically consists of three stages:

- **Unsupervised Pretraining.** The model learns general language patterns from a massive corpus of text. The objective is to predict tokens in a self-supervised manner, such as through autoregressive modeling (e.g., GPT [Radford \(2018\)](#)) or masked language modeling (e.g., BERT [Devlin et al. \(2019\)](#)). The resulting model, known as the base model, has broad linguistic knowledge but lacks task-specific fine-tuning.
- **Supervised Fine-Tuning (SFT).** The base model is further trained on a curated dataset with labeled examples. This process aligns the model with specific tasks, improving its ability to generate contextually appropriate responses. SFT helps refine the model’s behavior but does not fully optimize it for human preferences.
- **Reinforcement Learning with Human Feedback (RLHF).** This final step fine-tunes the model using reinforcement learning, guided by human feedback. Typically, a reward model is trained to rank responses based on quality, and the LLM is optimized using a technique called Proximal Policy Optimization (PPO; [Schulman et al. \(2017\)](#)). RLHF helps improve response alignment with human intent, making the model more helpful and less prone to undesirable behaviors.

The pretraining phase is the most computationally expensive, time-consuming, and financially demanding stage. It requires training on massive datasets using thousands of GPUs over weeks or even months. The large amount of data and model parameters leads to significant energy consumption and infrastructure costs, often amounting to ten million dollars or even more.

A recent LLM, DeepSeek-R1 ([Guo et al. \(2025\)](#)), released on January 20, 2025, has attracted much attention. According to its creators, DeepSeek-R1 was trained in just two months at a cost of 6 million dollars—significantly cheaper than OpenAI’s o1—while achieving comparable performance. Concerns over potential misallocation of investment in AI training have even contributed to a dip in NVIDIA’s stock price.

According to the original paper [Guo et al. \(2025\)](#), DeepSeek-R1 uses DeepSeek-V3-Base as its base model. A crucial question is whether the stated 6 million dollars includes the cost of pretraining. If not, the authors have primarily introduced a new pipeline to enhance LLMs’ reasoning capabilities, and many may have misunderstood this as a claim of overall cost efficiency. The real novelty might lie in their pipeline rather than in making the entire training process exceptionally cost-effective.

Market sometimes overacts, and social media often exaggerates. We still need some time to uncover the full picture.

6 Statistics in the New Era

Can Deep Learning Beat Statistics? Based on my personal experience, deep learning methods are powerful, but not as powerful as many people claim to be. Deep learning does not easily outperform traditional statistical methods in many real-world scenarios. For example, during my research internship at UM, I found that the performance of a vanilla linear regression model was comparable to that of a neural network when predicting the energy released by a solar flare event. More surprisingly, a standard logistic regression model slightly outperformed the neural network when classifying flare events into pre-defined categories. From a statistical perspective, neural networks are essentially over-parameterized models. When the sample size is not sufficiently large, we cannot expect such complex models to beat classical statistical models.

Challenges and Opportunities. Nowadays, many data are collected through automated processes (automatically, rather than manually) and across different generations of technology²⁴. As a result, the quality of the data is often low and measurement errors are inevitable, which calls for effective statistical inference across different scientific contexts. In addition, there is often a gap between the findings published in statistics journals and the concerns of practitioners. In the era of Big Data and AI, bridging this gap becomes even more essential. Future statisticians should have the ability to conduct end-to-end research, from data collection and preprocessing to model development, goodness-of-fit assessments, performance evaluation, and ultimately, communicating findings in ways that are impactful for real-world applications.

Statistics Ph.D. Students in the New Era. With the huge empirical success of large language models, diffusion models, and many other advancements in the era of AI, it would be unwise to ignore them. In my view, future statistics Ph.D. students should actively engage with other communities, including data science, machine learning, and computer science. We should join their seminars, exchange ideas during discussions, stay informed about their recent advancements, and contribute to their journals. Here is a quote from [Thomas Henry Huxley](#):

“Try to learn something about everything and everything about something.”

²⁴For example, during my research internship at UM, I learned that the detection algorithms for flare events have evolved over the years.

References

- Edo M Airoldi, David Blei, Stephen Fienberg, and Eric Xing. Mixed membership stochastic blockmodels. *Advances in neural information processing systems*, 21, 2008.
- Sanjeev Arora, Rong Ge, and Ankur Moitra. Learning topic models—going beyond svd. In *2012 IEEE 53rd annual symposium on foundations of computer science*, pages 1–10. IEEE, 2012.
- Leonard E Baum and Ted Petrie. Statistical inference for probabilistic functions of finite state markov chains. *The annals of mathematical statistics*, 37(6):1554–1563, 1966.
- Leonard E Baum, Ted Petrie, George Soules, and Norman Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *The annals of mathematical statistics*, 41(1):164–171, 1970.
- Julian Besag. Spatial interaction and the statistical analysis of lattice systems. *Journal of the Royal Statistical Society: Series B (Methodological)*, 36(2):192–225, 1974.
- Peter J Bickel and Elizaveta Levina. Covariance regularization by thresholding. *The Annals of Statistics*, 2008a.
- Peter J Bickel and Elizaveta Levina. Regularized estimation of large covariance matrices. *The Annals of Statistics*, 2008b.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- George EP Box and Norman R Draper. *Empirical model-building and response surfaces*. John Wiley & Sons, 1987.
- Stanley H Chan. Tutorial on diffusion models for imaging and vision. *arXiv preprint arXiv:2403.18103*, 2024.
- Minshuo Chen, Kaixuan Huang, Tuo Zhao, and Mengdi Wang. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In *International Conference on Machine Learning*, pages 4672–4712. PMLR, 2023.
- Corinna Cortes. Support-vector networks. *Machine Learning*, 1995.
- David R Cox. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2):187–202, 1972.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977.

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019. URL <https://api.semanticscholar.org/CorpusID:52967399>.
- David Donoho and Victoria Stodden. When does non-negative matrix factorization give a correct decomposition into parts? *Advances in neural information processing systems*, 16, 2003.
- Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011.
- Jianqing Fan, Yingying Fan, and Jinchi Lv. High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*, 147(1):186–197, 2008.
- Jianqing Fan, Yuan Liao, and Martina Mincheva. Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 75(4):603–680, 2013.
- Jianqing Fan, Runze Li, Cun-Hui Zhang, and Hui Zou. *Statistical foundations of data science*. Chapman and Hall/CRC, 2020.
- Howard A Garcia. Temperature and emission measure from goes soft x-ray measurements. *Solar Physics*, 154:275–308, 1994.
- Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. *Neural computation*, 4(1):1–58, 1992.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- T Hofmann. Probabilistic latent semantic indexing. In *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, 1999.
- Paul W Holland, Kathryn Blackmond Laskey, and Samuel Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- Pengsheng Ji, Jiashun Jin, Zheng Tracy Ke, and Wanshan Li. Co-citation and co-authorship networks of statisticians. *Journal of Business & Economic Statistics*, 40(2):469–485, 2022.

- Zhenbang Jiao, Hu Sun, Xiantong Wang, Ward Manchester, Tamas Gombosi, Alfred Hero, and Yang Chen. Solar flare intensity prediction with machine learning models. *Space weather*, 18(7):e2020SW002440, 2020.
- Jiashun Jin. Fast community detection by score. *The Annals of Statistics*, 2015.
- Jiashun Jin, Zheng Tracy Ke, and Shengming Luo. Estimating network memberships by simplex vertex hunting. *arXiv preprint arXiv:1708.07852*, 12, 2017.
- Jiashun Jin, Zheng Tracy Ke, and Shengming Luo. Mixed membership estimation for social networks. *Journal of Econometrics*, 239(2):105369, 2024.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37:183–233, 1999.
- Brian Karrer and Mark EJ Newman. Stochastic blockmodels and community structure in networks. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 83(1):016107, 2011.
- Zheng Tracy Ke and Minzhe Wang. Using svd for topic modeling. *Journal of the American Statistical Association*, 119(545):434–449, 2024.
- Zheng Tracy Ke, Pengsheng Ji, Jiashun Jin, and Wanshan Li. Recent advances in text analysis. *Annual Review of Statistics and Its Application*, 11, 2023.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Diederik P Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.
- Olivier Ledoit and Michael Wolf. A well-conditioned estimator for large-dimensional covariance matrices. *Journal of multivariate analysis*, 88(2):365–411, 2004.
- Tianxi Li, Elizaveta Levina, and Ji Zhu. Prediction models for network-linked data. *The Annals of Applied Statistics*, 2019.
- Jun S Liu. *Monte Carlo strategies in scientific computing*, volume 10. Springer, 2001.
- Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification from labeled and unlabeled documents using em. *Machine learning*, 39:103–134, 2000.
- Ian Porteous, David Newman, Alexander Ihler, Arthur Asuncion, Padhraic Smyth, and Max Welling. Fast collapsed gibbs sampling for latent dirichlet allocation. In *Proceedings of the*

- 14th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 569–577, 2008.
- Lawrence Rabiner and Biinghwang Juang. An introduction to hidden markov models. *ieee assp magazine*, 3(1):4–16, 1986.
- Alec Radford. Improving language understanding by generative pre-training. 2018.
- Adam J Rothman, Elizaveta Levina, and Ji Zhu. Generalized thresholding of large covariance matrices. *Journal of the American Statistical Association*, 104(485):177–186, 2009.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- Masashi Sugiyama and Motoaki Kawanabe. *Machine learning in non-stationary environments: Introduction to covariate shift adaptation*. MIT press, 2012.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- John W Tukey. Which part of the sample contains the information? *Proceedings of the National Academy of Sciences*, 53(1):127–134, 1965.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Jiangzhou Wang, Jingfei Zhang, Binghui Liu, Ji Zhu, and Jianhua Guo. Fast network community detection with profile-pseudo likelihood methods. *Journal of the American Statistical Association*, 118(542):1359–1372, 2023.
- Linshanshan Wang, Xuan Wang, Katherine P Liao, and Tianxi Cai. Semisupervised transfer learning for evaluation of model classification performance. *Biometrics*, 80(1):ujae002, 2024.

Frank Woodcock. The evaluation of yes/no forecasts for scientific and administrative purposes. *Monthly Weather Review*, 104(10):1209–1214, 1976.

Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.